

The allometry of cellular DNA and ribosomal gene content among microbes and its use for the assessment of microbiome community structure

Luis Gonzalez de Salceda

Arizona State University <https://orcid.org/0000-0002-9048-6264>

Ferran Garcia-Pichel (✉ ferran@asu.edu)

Arizona State University <https://orcid.org/0000-0003-1383-1981>

Research

Keywords: Bacteria, archaea, protists, fungi, genomes, ploidy, microbiomes, ribosomal genes

Posted Date: February 11th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-208798/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Microbiome on August 17th, 2021. See the published version at <https://doi.org/10.1186/s40168-021-01111-z>.

1 **The allometry of cellular DNA and ribosomal gene content among microbes and its use**
2 **for the assessment of microbiome community structure**

3 Luis González de Salceda and Ferran Garcia-Pichel

4 Center for Fundamental and Applied Microbiomics and School of Life Sciences, Arizona State
5 University

6 Email addresses of the authors: imgonz21@asu.edu and ferran@asu.edu

7 Corresponding author: Ferran Garcia-Pichel

8

9

10

11

12

13

14

15

16

17

18

19

20

21 **Abstract**

22 **Background.** The determination of taxon-specific composition of microbiomes by combining
23 high-throughput sequencing of ribosomal genes with phyloinformatic analyses has become
24 routine in microbiology and allied sciences. Systematic biases to this approach based on the
25 demonstrable variability of ribosomal operon copy number per genome were recognized early.
26 The more recent realization that polyploidy is probably the norm, rather than the exception,
27 among microbes from all domains of life, points to an even larger source bias.

28 **Results.** We found that the number of 16S or 18S RNA genes per cell, a combined result of the
29 number of RNA gene loci per genome and ploidy level, follows an allometric power law of cell
30 volume with an exponent of $2/3$ across 6 orders of magnitude in small subunit copy number per
31 cell and 9 orders of magnitude in cell size. This stands in contrast to cell DNA content, which
32 follows a power law with an exponent of $3/4$.

33 **Conclusion.** In practical terms, that relationship allows for a single, simple correction for
34 variations in both copy number per genome and ploidy level in ribosomal gene analyses of taxa-
35 specific abundance. In biological terms, it points to the uniqueness of ribosomal gene content
36 among microbial properties that scale with size.

37

38

39 **Key Words.** Bacteria, archaea, protists, fungi, genomes, ploidy, microbiomes, ribosomal genes.

40

41

42

43 **Background**

44 The rRNA gene approach to microbiome analyses, either based on amplicon or metagenomic
45 sequencing, relies on the tacit assumption that the counts of these marker gene translate into a
46 robust measure or proxy for microbial abundance. However, this assumption is often violated.
47 Sources of error in gene abundance determinations can come from analytical procedures such
48 as DNA extraction, PCR amplification and sequencing itself (1). But likely as important,
49 systematic biases can be caused by the varying abundance of ribosomal genes in the genomes
50 of microbes (2). The concern is evident in the dedicated databases that document the variability
51 in ribosomal copy number per genome (R_g) among microbes (3). Interestingly, R_g seems to
52 correlate with a microbe's life history traits, where fast growth is associated with higher values
53 (4-6). There is also evidence for a certain degree of conservation in R_g within bacterial
54 phylogenetic clades (7). On this basis, bioinformatic tools have been developed to automatically
55 correct ribosomal surveys for R_g (8). The phylogenetic conservation of R_g , however, seems only
56 conspicuous among closely related microbes (9), and can explain only some 10 % of its
57 variability in complex, diverse communities (10). In some eukaryotes like *Saccharomyces*
58 *cerevisiae* R_g is unstable and can vary widely among strains or individuals (11). Importantly,
59 such corrections would only lead us to a description of community composition in terms of
60 relative abundance of taxon-specific genome copies. But more useful metrics in microbiome
61 community composition analyses are either cell number (7) (i.e, individuals) or biomass
62 contributions by each taxon. Given the close to 9 orders of magnitude spanned by microbial cell
63 biomass, it can be argued that taxon-specific biomass rather than cell number would be a better
64 descriptor of a taxon's contribution to a community. However, there are still instances where
65 number of cells would be preferred (for example, to gauge dispersal potential, culturability, or
66 susceptibility to deleterious agents like predators or toxicants). In any case, to translate genome
67 numbers to cell numbers one needs to take into account the level of ploidy, P , the number of

68 copies of the genome present in a cell, where the number of ribosomal operons per cell, R_c is
69 the product PR_g . Surprisingly, P is not typically taken into account, perhaps under the
70 assumption that most microbes, like *Escherichia coli*, are monoploid (12)(13). And yet, in
71 bacteria and archaea, P varies far more than R_g (12, 14), and most species examined are oligo-
72 or polyploid, with some containing in excess of 200 genomes copies per cell (15). If one
73 includes unicellular eukaryotes, the variation can be 4 orders of magnitude (16). Clearly, ploidy
74 constitutes a very important source of bias for community counts in itself (17), affecting
75 estimates from both amplicon sequencing and shot-gun metagenomics. The variable nature of
76 P could potentially either compound or diminish the effect of R_g variability in determining a cell's
77 R_c , as it is not known whether P and R_g correlate or vary independently among species; a high
78 P could be associated with low R_g , and vice-versa. Studies on marine protists intended to
79 estimate biomass from 18S counts have shown that R_c correlated linearly with cell volume (V_c)
80 (18) or cell length (19) when plotted on double log scales, indicating an R_c dependence on size.
81 Here, we posited that perhaps there is constancy among microbes in the need for ribosomal
82 gene content in relation to their cell biomass. In other words, microbial species would be under
83 selection to contain a sufficient but not excessive R_c to support the production of their typical cell
84 biomass, B_c , so that R_c would be proportional to B_c . Assuming cell density to be invariant
85 (around 1.008 g ml^{-1}) (20), R_c would also be proportional to cell volume (V_c).

86

87 **Methods**

88 **Dataset.** Values for all parameters were gathered or derived from the literature. In place of B_c ,
89 we used cellular volume, V_c , assuming cellular density to be constant (around 1.008 g ml^{-1} (20)).
90 Cell volumes were either taken from reported direct determinations or derived from literature
91 photomicrographs assuming simple formulae for a variety of fitting three-dimensional shapes

92 (i.e sphere, cylinder) or combinations thereof as given in Table S1 (see Additional file 1). When
93 a range of volume values was available, we used the average. For R_g , we used values given in
94 rrnDB (3) for the same species or strain. If they were not available, we used literature values or
95 determined it by examination of the strain's publicly available genome through BLAST. Ploidy
96 was either taken directly from reported values or estimated if cellular DNA content *and* genome
97 size were known. If P was variable within a species or strain, we used the average level of the
98 range given. R_c values were then derived as the product of P and R_g , although for many protists,
99 R_c was taken directly from experimentally determined values. The annotated input data are
100 gathered in Table S1 (see Additional file 1). The limiting factor to the size of the database was
101 the availability of P determinations, which are quite uncommon. In all, we could analyze 107
102 cases.

103 **Statistics.** Power fits of data were run in Excel as linear regressions of the ln-transformed data
104 pairs using a least-squares model. Statistics are given in Table S2 (see Additional file 2). To test
105 the significance of exponent differences in two separate datasets, we used T-tests for the
106 slopes of the linear fits.

107 **Estimation of taxon-specific cell numbers and biovolumes from 16S rRNA counts.** In a
108 dataset of rRNA gene taxon-specific frequencies, F_r , assigned to i taxa whose cell volumes,
109 $V_c(i)$, are known, one can directly estimate $R_c(i)$ from Equation 1 (see Results section). The
110 relative contribution to number of cells by taxon i , $F_c(i)$, is computed as:

$$111 \quad F_c(i) = \frac{F_r(i)}{R_c(i) \sum R_c(i)}$$

112 And the relative contribution to biovolume, $F_v(i)$, as:

$$113 \quad F_v(i) = \frac{F_c(i) V_c(i)}{\sum F_c(i) V_c(i)}$$

114 If a determination of the absolute abundance of the total copies of the ribosomal gene for all
115 taxa considered in the sample of origin, R_s , is available (from qPCR, for example, in units of
116 copies per mass, volume or surface sampled), then absolute taxon-specific assignments $R(i)$,
117 can be obtained as the product $F_r(i)R_s(i)$. From $R(i)$ one can derive absolute values for cells
118 $C(i)$ and biovolume $V(i)$ attributable to each taxon: $C(i) = R(i) / R_c(i)$, and $V(i) = C(i)V_c(i)$. The
119 sums $\sum C(i)$ and $\sum V(i)$, estimate the absolute number of cells or biovolume (in μm^3),
120 respectively, of the entire set of taxa under consideration.

121 An alternative to using $V_c(i)$, if those are not exactly known, is to assign rough discrete size
122 ranges to taxa, and to use mean V_c (and R_c) values of the range's maximum and minimum. We
123 found it advisable to set variable-width size ranges in such a way that within-range variation in
124 resulting R_c values was kept moderate. We used the following cell diameter ranges (in μm): 0.2-
125 0.3, 0.3-0.4, 0.4-0.6, 0.6-0.9, 0.9-1.2, 1.2-1.5, 1.5-2.1, 2.1-2.9, 2.9-4.1, 4.1-5.8, 5.8-8.2, 8.2-11.6,
126 11.6-16.4. This set provides within-range variation in R_c of less than 8% in all cases, which is
127 smaller than the uncertainty of our estimates for the normalization constant in Eq. 1 of Results.

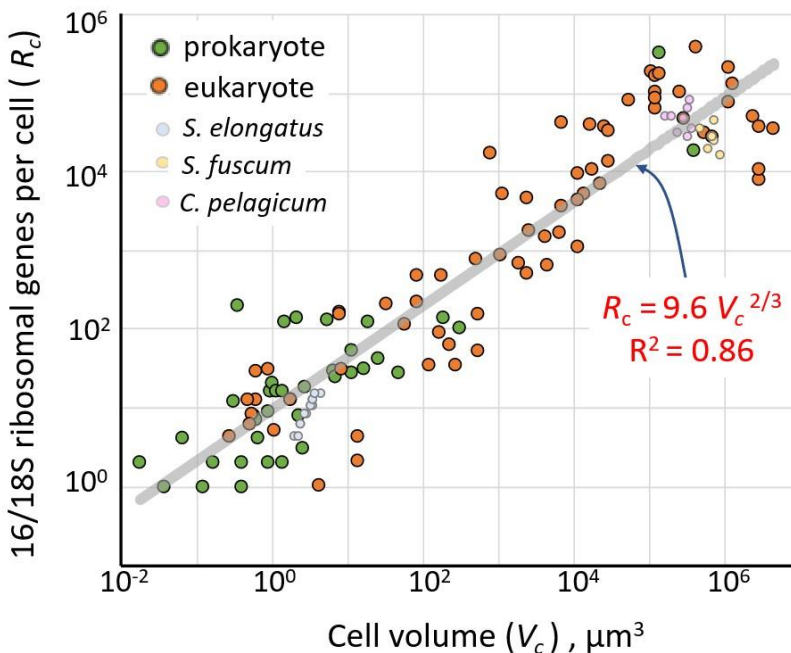
128

129 **Results**

130 Traits that span orders of magnitude are best evaluated as double logarithmic plots, which can
131 be analyzed by power function fits. In this approach, the hypothesis of proportionality between
132 V_c and R_c we posed should have resulted in a power function fit with an exponent close to unity.
133 Our analysis (Fig. 1) readily dispelled that contention. The fit instead revealed that R_c follows
134 well ($R^2 = 0.86$) a power function of V_c with an exponent significantly lower than unity, and
135 indistinguishable from 2/3 (0.66 ± 0.03 ; \pm SE) across nine orders of magnitude in cell volume.
136 For volumes expressed in μm^3 ,

$$137 \quad R_c = 9.58 V_c^{0.66} \cong 9.58 V_c^{2/3} \quad \text{[Equation 1],}$$

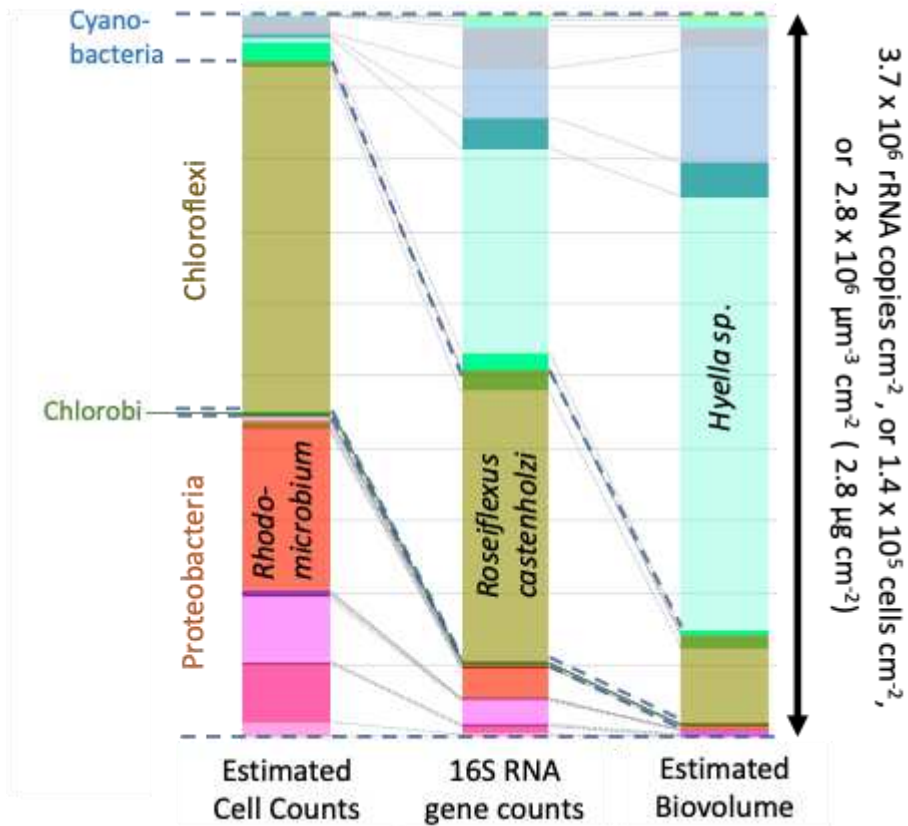
138 where 9.58 ± 1.21 is the estimated normalization constant. One could envision that the scaling
 139 relationship may have been artifactually distorted at the low range of R_c , since it cannot
 140 physiologically take values < 1 . But a reanalysis of the dataset excluding data pairs with $R_c \leq 2$
 141 did not change the fit significantly in exponent or normalization constant (see Additional file 2).
 142 We also tested the hypothesis that exponents for a fit of data pairs from eukaryotes (exponent =
 143 0.72 ± 0.05) vs. prokaryotes (0.62 ± 0.05) could be different, but this did not find strong
 144 statistical support in a T-test comparison ($p = 0.20$).



145 **Figure 1.** Relationship between cellular ribosomal gene content (R_c) and cell volume (V_c) in
 146 microbes ($n = 107$), plotted as a log/log graph. The grey line is a power fit with the equation
 147 displayed in red type (fit statistics are in Table S2, Additional file 2). Datapoints belonging to
 148 eukaryotes are in orange, those for prokaryotes in green. For three species, we plotted datasets
 to highlight intraspecies variability: *Synechococcus elongatus* (light blue symbols) (28), *Colozoum*
pelagicum (light purple) (19), and *Sphaerozoum fuscum* (19) (light yellow).

149 Thus, R_c scales generally not with the volume but with the surface area of a microbial cell, which
 150 for the purpose of this study means that the bias associated with ribosomal gene counts will be
 151 size-dependent regardless of our choice of abundance estimator. Ribosomal counts will

152 overestimate large-celled microbes over small-celled ones if one is interested in number of cells,
153 the bias increasing with the square of linear cell size (Eq. 2), a prediction that finds experimental
154 support for specific cases in the literature (21). In terms of biomass, ribosomal counts will
155 underestimate the contribution of large microorganisms, the bias increasing with the $2/3$ power
156 of cellular biovolume (Eq. 1). Whichever the desired measure of abundance, however, the
157 explicit relationship in Fig. 1 provides a means for bias correction in tallies of ribosomal genes,
158 as long as cell biovolume is known from ancillary data for the taxa detected in the microbiome of
159 interest. The correction requires knowledge of neither P nor R_g .



160

Figure 2. Estimation of microbial community structure based on experimental ribosomal counts (central column), estimated cell number (left column) and estimated biovolume (right column) in a single, exemplary dataset using allometric corrections based on Eq. 1. The dataset is from Roush et al. (2020) (22), and includes the subset of taxonomically assignable phototrophic bacteria from an endolithic microbiome on coastal marine carbonate rocks. Only three exemplary phototrophs are labeled, but full, taxonomically explicit distributional data are in Table S3 (see Additional file 3). For ease of comparison, results are graphically presented as relative frequencies, but absolute scales of areal abundance are indicated on the arrow to the right.

161

162 A procedural explanation is given under Methods, and we provide an example application in Fig.
 163 2 using a dataset of phototrophic bacteria from endolithic microbiomes within intertidal hard
 164 carbonate rocks (22), responsible for their micritization and bioerosion (23), and useful here
 165 because typical cell volumes could be assigned to all taxa. The differential outcomes are

166 obvious: 16S rRNA counts of large-celled cyanobacterial genera severely underestimate their
167 contribution to biomass but overestimate their contribution in terms of number of cells (see for
168 example, *Hyella* sp.). The opposite is true for alphaproteobacterial phototrophs (see for example
169 *Rhodospirillum rubrum* sp.), most of which are small-celled (24). The distortion is less intense for the
170 Chloroflexi, with intermediate cell size (see *Roseiflexus castenholzii*, for example).

171 We have presented the issue of bias having in mind relative abundance tallies of microbiome
172 members, but proportional tallies have methodological constraints in themselves, because the
173 individual proportions must add up to 1, and thus the relative abundances of taxa are
174 necessarily not independent of each other. There is clear evidence of severely diverging
175 analytical outcomes when both relative and absolute abundance are compared in the same
176 datasets (25, 26). Commonly, relative proportions or taxa-specific ribosomal copies are
177 converted to absolute abundances with parallel quantification of rRNA gene copies by qPCR,
178 either total copies in the community analyzed or those of particular taxa (16). We note here that,
179 in view of our results, the latter would require allometric correction, whereas the former would
180 not (as done in the dataset presented in Fig. 2) and is thus a preferable approach. However, we
181 also note that the total number of ribosomal gene copies in a sample is not a good absolute
182 measure of the combined microbial biomass or number of cells present for comparisons among
183 samples, as it will be dependent on their inherent cell-size distribution. Hence, comparisons
184 among samples will only be meaningful if carried out after conversion to biomass or cell
185 numbers, unless the microbial composition of the samples is unchanged.

186

187 **Discussion**

188 The procedure outlined here requires knowledge of morphological metadata in addition to
189 sequencing counts for each taxon. Unfortunately, cell volume data are not readily available for

190 many taxa, at least in a compiled format, and requires intensive literature searches. In its
191 absence, and as an approximation, using a few discrete cell-size classes instead of exact
192 values yields useful corrected distributions (see Additional file 4). Yet, an effort to bring microbial
193 size data into a consolidated platform would be desirable in that it would enable the processing
194 of large datasets in an automated, more manageable way.

195 An additional factor to take into account is the substantial data spread around the fit leading to
196 Eq. 1, which can limit the precision of the correction. An expanded dataset should improve
197 predictive accuracy and perhaps even precision, but some inherent limitations are also at play.
198 P can vary in a single strain with cell cycle (15) and growth conditions (27). We have included
199 the range of intraspecies variability on the V_c/R_c space in Fig. 1 for the cases a single strain of
200 *Synechococcus elongatus*, and single cells from natural populations of *Sphaerozoum fuscum*
201 and *Colozoum pelagicum*. They suggest that a significant proportion of spread can be attributed
202 to biological intraspecies variability, tempering the prospects for improvement with eventually
203 extended datasets. Studies on *Synechococcus elongatus* (28, 29) (30) and *Saccharomyces*
204 *cerevisiae* (31) point to a regulatory interdependency of P with cell size. Additionally, because
205 the data used here were arrived at through several approaches, a dedicated survey based on
206 more consistent analytical procedures may result in tighter fits. Finally, part of the variability
207 detected may have been due to neglecting contributions of organelle ribosomal genes in
208 protists. This is expected to be negligible for large-celled eukaryotes, but perhaps not so much
209 for the smallest of them, in which organelles take up a larger portion of their cell volume. Indeed,
210 some of these pico-eukaryotes contribute disproportionately (by defect in R_c) to the regression's
211 sum of squares and may have contributed to the somewhat higher exponent in the eukaryote-
212 only fit (Additional file 2). In support of this contention, a re-analysis excluding eukaryotes with
213 $V_c < 20 \mu\text{m}^3$ yields an exponent (0.66 ± 0.06 ; $R^2 = 0.68$), more in line with that of Eq. 1, showing

214 no trace of statistical difference (T-test, $p = 0.68$) with that of the prokaryote-only fit (Additional
215 file 2).

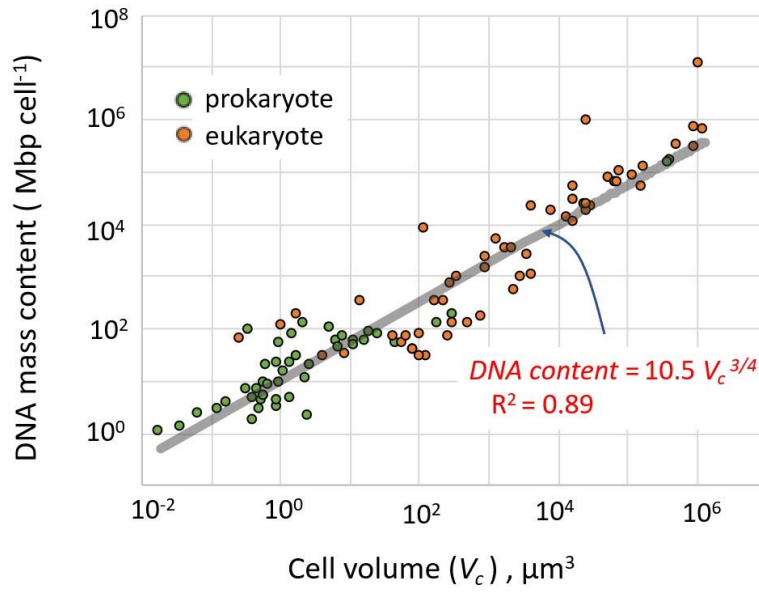
216 The preceding discussion on uncertainty in the correction approach should not be taken as
217 grounds for inaction, given that the range of variation in V_c far exceeds that traceable to
218 deviations from the fit, not only among microbes at large, but also in specific settings, and the
219 spectra of microbial size distribution in microbiomes seems to be dynamic. For example, the
220 range of V_c of typical bacterioplankton (excluding phototrophs) in seawater spans 3 orders of
221 magnitude and its spectrum can be modified significantly by factors like grazing (32).
222 Considering photosynthetic plankton would likely add another 4-5 orders of magnitude in V_c ,
223 and the size spectrum of this group is also affected by environmental parameters (33). In the
224 human gut microbiome, our initial assessments show that microbiome typical bacteria span over
225 at least 4 orders of magnitude in volume.

226 Beyond the pragmatic uses for community composition corrections, we see it as unlikely that the
227 apparent scaling relationship with cell surface area has no biological meaning. It is tempting to
228 speculate that R_c scales with size to maintain an increasing protein content need. Indeed,
229 protein content scales as a function of cell volume with a similar exponent of 0.70 ± 0.06 ($R^2 =$
230 0.87 ; 95% CI: 0.64 - 0.75) (34). Because the CI of the exponent for protein content per cell and
231 that for R_c in the fit of Fig. 1 [0.72 and 0.61; see Additional file 2] overlap, the possibility of a
232 connection to cellular need for proteins cannot be rejected solely on this basis. Indeed, in
233 *Synechococcus elongatus* in the laboratory, protein content and P (hence also R_c) strongly co-
234 vary with cell volume (30).

235 Alternatively, and perhaps more trivially, the scaling relationship of R_c with V_c may simply be a
236 reflection of the size scaling of DNA content per cell. In other words, ribosomal genes would
237 follow the trends of DNA content as a whole, just like any other universal gene would. The

238 allometric relationship between DNA content and cell size, however, has not been addressed in
239 the literature or has been addressed incorrectly by neglecting ploidy (34, 35). We know that
240 genome size scales among bacteria with reported exponents between 0.21 ($R^2 = 0.60$) (34) and
241 0.35 ($R^2 = 0.45$) (36). In our dataset, which includes eukaryotes, it does so with an exponent of
242 0.18 ($R^2 = 0.34$; see Additional file 5). Even when these coefficients of correlation are rather
243 poor, genome size clearly increases much more weakly with V_c than does R_c . Again, this does
244 not take into account P variations to yield actual DNA content per cell; it is the size of one copy
245 of the genome. A portion of our dataset can be used to explore the scaling of DNA content per
246 cell for prokaryotes ($n = 60$). To this subset we can add the measured or slightly derived values
247 reported by Shuter et al (35) ($n=39$), excluding those that relied on assumptions of monoploidy.
248 This combined set yields a power scaling fit with $R^2 = 0.89$ and estimated exponent of $\frac{3}{4}$ ($0.75 \pm$
249 0.03 ; Fig 3).

250



251

Figure 3. DNA content scales with cell volume as a power function with an exponent of $3/4$. Entries are from a subset of those in Table S1 ($n = 60$, see Additional file 1), and determinations by Shuter et al (35) ($n=39$). Orange points belong to eukaryotic microbes and green points to prokaryotes. Full statistics for the fit (in red type) are given in Table S2 (Additional file 2).

252

253 The difference in scaling exponent between genome size and cell DNA content (0.18-035 vs.
 254 0.75) gauges the importance of P . In fact, in our dataset, P seems to scale with V_c as a power
 255 law with an exponent of 0.54 ($R^2 = 0.69$; Additional file 6). This is consistent with the fact that
 256 the product of genome size and ploidy yields the cell DNA content, as the exponents of the
 257 multipliers (0.18 and 0.54, respectively) roughly add up to the estimated exponent of the product
 258 (0.75). That the exponents for DNA content per cell ($3/4$) and R_c ($2/3$) are significantly different
 259 (T-test, $p = 0.02$), speaks for respective mechanistic drivers that are fundamentally decoupled.
 260 In fact, most known allometric laws found in nature scale with exponents that are simple
 261 multiples of $1/4$ (37). It would seem that ribosomal genes are, in that sense, unique.

262

263

264 **Conclusions**

265 The results presented here uncover surprising basic rules on the composition of microbes, rules
266 that ties them all together, and that far from being self-evident, pose an intellectual challenge to
267 elucidate. In practical terms, this discovery also provides a rather simple approach to deal with
268 biases affecting the use of current omics methodologies for the assessment of microbiome
269 composition, which, given their extensive use in many areas of microbiology and allied
270 sciences, has a large potential for applicability.

271 **Additional Files**

272 **Additional file 1: Table S1.** Taxon-specific values for primary variables (cell volume, ribosomal
273 gene copies per cell, cellular DNA content and ploidy) as well as source variables (cell shape,
274 cell axial dimension, ribosomal gene copies per genome and genome size) as used in the
275 analyses presented in Figure 1. (XLSX 22 kb)

276 **Additional file 2: Table S2.** Statistics and estimated parameters for power fits against V_c .
277 (DOCX 14 kb)

278 **Additional file 3: Table S3.** Explicit dataset used in Figure 2. Original 16S rRNA gene amplicon
279 sequencing data, taxonomic assignments, and qPCR 16S rRNA gene quantifications are from
280 Roush et al. (2020). Estimations of taxon-specific cell number and taxon-specific biovolume
281 according to Materials and Methods. (XLSX 14 kb)

282 **Additional file 4: Figure S1.** Differences in allometric estimation of microbial community
283 structure as cell number or biovolume from 16S rRNA gene counts in the dataset of Fig. 2 by
284 either assigning measured cell volume values to taxa or by assigning taxa to a set of discrete

285 size ranges. Left: stack bar graphs for relative proportions of taxa. Right: frequency histograms
286 for taxa-specific percentual differences between the two approaches. (DOCX 66 kb)

287 **Additional file 5: Figure S2.** Relationship between genome size and cell volume (V_c) in
288 microbes ($n = 56$), plotted as a log/log graph. The grey line is a power fit with the equation
289 displayed in red type (fit statistics are in Suppl. Table 2). Datapoints belonging to eukaryotes are
290 in orange, those for prokaryotes in green. (DOCX 36 kb)

291 **Additional file 6: Figure S3.** Relationship between ploidy (P) and cell volume (V_c) in microbes
292 ($n = 56$), plotted as a log/log graph. The grey line is a power fit with the equation displayed in
293 red type (fit statistics are in Suppl. Table 2). Datapoints belonging to eukaryotes are in orange,
294 those for prokaryotes in green. (DOCX 25 kb)

295

296 **List of abbreviations**

297 R_g : Ribosomal gene copy number per genome

298 R_c : Ribosomal gene copy number per cell

299 R_s : Total copies of the ribosomal gene for all taxa considered in a sample

300 V_c : Cell volume

301 B_c : Cell biomass

302 D_c : Cell diameter

303 P : Ploidy

304

305 **Declarations**

306

307 **Ethics approval and consent to participate.** Not applicable

308

309 **Consent for publication.** Not applicable

310

311 **Availability of data and materials.** All data generated or analyzed during this study are

312 included in this published article and its supplementary information files

313

314 **Competing interests.** The authors declare that they have no competing interests

315

316 **Funding** This work was supported in part by NSF grant #1449501

317

318 **Authors' contributions.** FGP conceived the idea. LGS carried out the database research. LGS

319 and FGP run the analyses and wrote the manuscript. All authors read and approved the final

320 manuscript.

321 **Acknowledgements.** We thank Susanne Neuer, Gillian Gile and Damien Finn for critically

322 reading the manuscript and D. Roush for walking us through the data used in the exemplary

323 corrections.

324

325 **References**

- 326 1. Brooks JP, *et al.* (2015) The truth about metagenomics: quantifying and counteracting
327 bias in 16S rRNA studies. *BMC microbiology* 15(1):1-14.
- 328 2. Lavrinienko A, Jernfors T, Koskimäki JJ, Pirttilä AM, & Watts PC (2020) Does
329 Intraspecific Variation in rDNA Copy Number Affect Analysis of Microbial Communities?
330 *Trends in microbiology*.
- 331 3. Stoddard SF, Smith BJ, Hein R, Roller BRK, & Schmidt TM (2014) rrnDB: improved tools
332 for interpreting rRNA gene abundance in bacteria and archaea and a new foundation for
333 future development. *Nucleic Acids Research* 43(D1):D593-D598.
- 334 4. Klappenbach JA, Dunbar JM, & Schmidt TM (2000) rRNA operon copy number reflects
335 ecological strategies of bacteria. *Applied and environmental microbiology* 66(4):1328-
336 1333.
- 337 5. Roller BR, Stoddard SF, & Schmidt TM (2016) Exploiting rRNA operon copy number to
338 investigate bacterial reproductive strategies. *Nature microbiology* 1(11):1-7.
- 339 6. Vieira-Silva S & Rocha EP (2010) The systemic imprint of growth and its uses in
340 ecological (meta) genomics. *PLoS Genet* 6(1):e1000808.
- 341 7. Kembel SW, Wu M, Eisen JA, & Green JL (2012) Incorporating 16S gene copy number
342 information improves estimates of microbial diversity and abundance. *PLoS Comput Biol*
343 8(10):e1002743.
- 344 8. Angly FE, *et al.* (2014) CopyRighter: a rapid tool for improving the accuracy of microbial
345 community profiles through lineage-specific gene copy number correction. *Microbiome*
346 2(1):1-13.
- 347 9. Lofgren LA, *et al.* (2019) Genome-based estimates of fungal rDNA copy number
348 variation across phylogenetic scales and ecological lifestyles. *Molecular ecology*
349 28(4):721-730.
- 350 10. Louca S, Doebeli M, & Parfrey LW (2018) Correcting for 16S rRNA gene copy numbers
351 in microbiome surveys remains an unsolved problem. *Microbiome* 6(1):41.

- 352 11. Kwan EX, Wang XS, Amemiya HM, Brewer BJ, & Raghuraman M (2016) rDNA Copy
353 number variants are frequent passenger mutations in *Saccharomyces cerevisiae*
354 deletion collections and de novo transformants. *G3: Genes, Genomes, Genetics*
355 6(9):2829-2838.
- 356 12. Pecoraro V, Zerulla K, Lange C, & Soppa J (2011) Quantification of ploidy in
357 proteobacteria revealed the existence of monoploid,(mero-) oligoploid and polyploid
358 species. *PloS one* 6(1):e16392.
- 359 13. Trun NJ (1998) Genome ploidy. *Bacterial Genomes*, (Springer), pp 95-102.
- 360 14. Soppa J (2014) Polyploidy in archaea and bacteria: about desiccation resistance, giant
361 cell size, long-term survival, enforcement by a eukaryotic host and additional aspects.
362 *Journal of molecular microbiology and biotechnology* 24(5-6):409-419.
- 363 15. Maldonado R, Jiménez J, & Casadesús J (1994) Changes of ploidy during the
364 *Azotobacter vinelandii* growth cycle. *Journal of bacteriology* 176(13):3911-3919.
- 365 16. Bonk F, Popp D, Harms H, & Centler F (2018) PCR-based quantification of taxa-specific
366 abundances in microbial communities: Quantifying and avoiding common pitfalls.
367 *Journal of microbiological methods* 153:139-147.
- 368 17. Soppa J (2017) Polyploidy and community structure. *Nature microbiology* 2(2):1-2.
- 369 18. Godhe A, *et al.* (2008) Quantification of diatom and dinoflagellate biomasses in coastal
370 marine seawater samples by real-time PCR. *Applied and environmental microbiology*
371 74(23):7174-7182.
- 372 19. Biard T, *et al.* (2017) Biogeography and diversity of Collodaria (Radiolaria) in the global
373 ocean. *The ISME Journal* 11(6):1331-1344.
- 374 20. Guerrero R, Pedrós-Alió C, Schmidt TM, & Mas J (1985) A survey of buoyant density of
375 microorganisms in pure cultures and natural samples. *Microbiologia (Madrid, Spain)* 1(1-
376 2):53-65.

- 377 21. Jasso-Selles DE, *et al.* (The Complete Protist Symbiont Communities of *Coptotermes*
378 *formosanus* and *Coptotermes gestroi*: Morphological and Molecular Characterization of
379 Five New Species. *Journal of Eukaryotic Microbiology* n/a(n/a).
- 380 22. Roush D & Garcia-Pichel F (2020) Succession and Colonization Dynamics of Endolithic
381 Phototrophs within Intertidal Carbonates. *Microorganisms* 8(2):214.
- 382 23. Chacón E, Berrendero E, & Pichel FG (2006) Biogeological signatures of microboring
383 cyanobacterial communities in marine carbonates from Cabo Rojo, Puerto Rico.
384 *Sedimentary Geology* 185(3-4):215-228.
- 385 24. Overmann J & Garcia-Pichel F (2006) The phototrophic way of life. *The prokaryotes*
386 2:32.
- 387 25. Props R, *et al.* (2017) Absolute quantification of microbial taxon abundances. *The ISME*
388 *Journal* 11(2):584-587.
- 389 26. Fernandes VM, *et al.* (2018) Exposure to predicted precipitation patterns decreases
390 population size and alters community structure of cyanobacteria in biological soil crusts
391 from the Chihuahuan Desert. *Environmental microbiology* 20(1):259-269.
- 392 27. Paranjape SS & Shashidhar R (2017) The ploidy of *Vibrio cholerae* is variable and is
393 influenced by growth phase and nutrient levels. *FEMS Microbiology Letters* 364(19).
- 394 28. Ohbayashi R, *et al.* (2019) Coordination of Polyploid Chromosome Replication with Cell
395 Size and Growth in a Cyanobacterium. *mBio* 10(2):e00510-00519.
- 396 29. Zheng X-y & O'Shea EK (2017) Cyanobacteria maintain constant protein concentration
397 despite genome copy-number variation. *Cell reports* 19(3):497-504.
- 398 30. Chen AH, Afonso B, Silver PA, & Savage DF (2012) Spatial and temporal organization
399 of chromosome duplication and segregation in the cyanobacterium *Synechococcus*
400 *elongatus* PCC 7942. *PLoS one* 7(10):e47837.
- 401 31. Mundkur BD (1953) Interphase nuclei and cell sizes in a polyploid series of
402 *Saccharomyces*. *Experientia* 9(10):373-374.

- 403 32. Andersson A, Larsson U, & Hagström Å (1986) Size-selective grazing by a
404 microflagellate on pelagic bacteria. *Marine Ecology Progress Series*:51-57.
- 405 33. Marañón E, *et al.* (2001) Patterns of phytoplankton size structure and productivity in
406 contrasting open-ocean environments. *Marine Ecology Progress Series* 216:43-56.
- 407 34. Kempes CP, Wang L, Amend JP, Doyle J, & Hoehler T (2016) Evolutionary tradeoffs in
408 cellular composition across diverse bacteria. *The ISME journal* 10(9):2145-2157.
- 409 35. Shuter BJ, Thomas J, Taylor WD, & Zimmerman AM (1983) Phenotypic correlates of
410 genomic DNA content in unicellular eukaryotes and other cells. *The American Naturalist*
411 122(1):26-44.
- 412 36. DeLong JP, Okie JG, Moses ME, Sibly RM, & Brown JH (2010) Shifts in metabolic
413 scaling, production, and efficiency across major evolutionary transitions of life.
414 *Proceedings of the National Academy of Sciences* 107(29):12941-12945.
- 415 37. West GB, Woodruff WH, & Brown JH (2002) Allometric scaling of metabolic rate from
416 molecules and mitochondria to cells and mammals. *Proceedings of the National*
417 *Academy of Sciences* 99(suppl 1):2473-2478.

418

419

420

421

422

423

424

425

426 **Figure titles and legends**

427 **Figure 1. Title:** Relationship between cellular ribosomal gene content (R_c) and cell volume (V_c)
428 in microbes.

429 **Legend.** The grey line is a power fit ($n = 107$) with the equation displayed in red type (fit
430 statistics are in Table S2, Additional file 2). Datapoints belonging to eukaryotes are in orange,
431 those for prokaryotes in green. For three species, we plotted datasets to highlight intraspecies
432 variability: *Synechococcus elongatus* (light blue symbols) (28), *Colozoum pelagicum* (light
433 purple) (19), and *Sphaerozoum fuscum* (19) (light yellow).

434 **Figure 2. Title:** Microbial community structure as determined by ribosomal counts, and as on
435 estimated on the basis of cell numbers and biovolume.

436 **Legend:** Estimation of microbial community structure based on experimental ribosomal counts
437 (central column), estimated cell number (left column) and estimated biovolume (right column) in
438 a single, exemplary dataset using allometric corrections based on Eq. 1. The dataset is from
439 Roush et al. (2020) (22), and includes the subset of taxonomically assignable phototrophic
440 bacteria from an endolithic microbiome on coastal marine carbonate rocks. Only three
441 exemplary phototrophs are labeled, but full, taxonomically explicit distributional data are in Table
442 S3 (see Additional file 3). For ease of comparison, results are graphically presented as relative
443 frequencies, but absolute scales of areal abundance are indicated on the arrow to the right.

444 **Figure 3. Title:** DNA content scales among microbes with cell volume as a power function with
445 an exponent of $\frac{3}{4}$.

446 **Legend:** Entries are from a subset of those in Table S1 ($n = 60$, see Additional file 1), and
447 determinations by Shuter et al (35) ($n=39$). Orange points belong to eukaryotic microbes and
448 green points to prokaryotes. Full statistics for the fit (in red type) are given in Table S2
449 (Additional file 2).

Figures

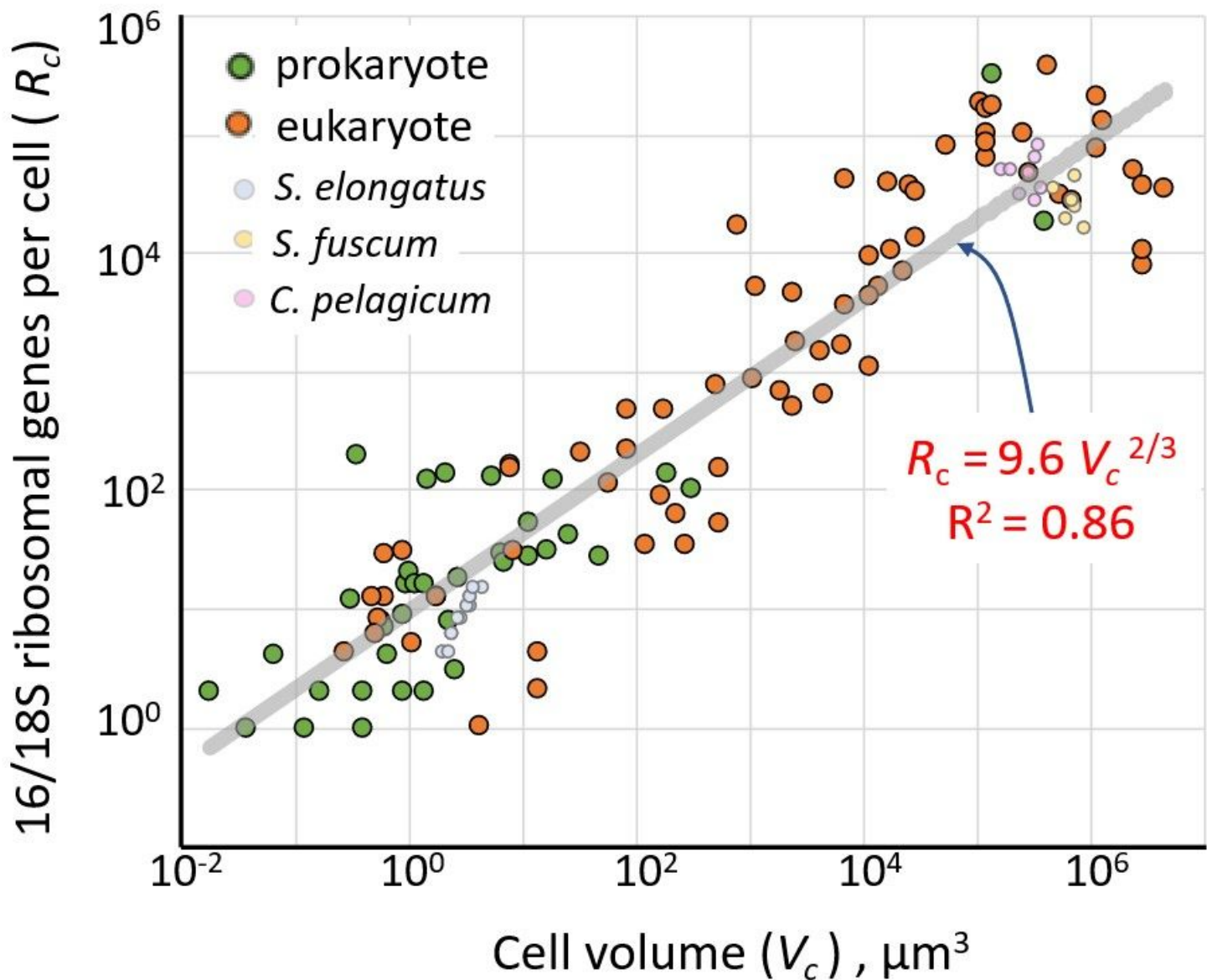


Figure 1

Relationship between cellular ribosomal gene content (R_c) and cell volume (V_c) in microbes. Legend. The grey line is a power fit ($n = 107$) with the equation displayed in red type (fit statistics are in Table S2, Additional file 2). Datapoints belonging to eukaryotes are in orange, those for prokaryotes in green. For three species, we plotted datasets to highlight intraspecies variability: *Synechococcus elongatus* (light blue symbols) (28), *Colozoum pelagicum* (light purple) (19), and *Sphaerozoum fuscum* (19) (light yellow).

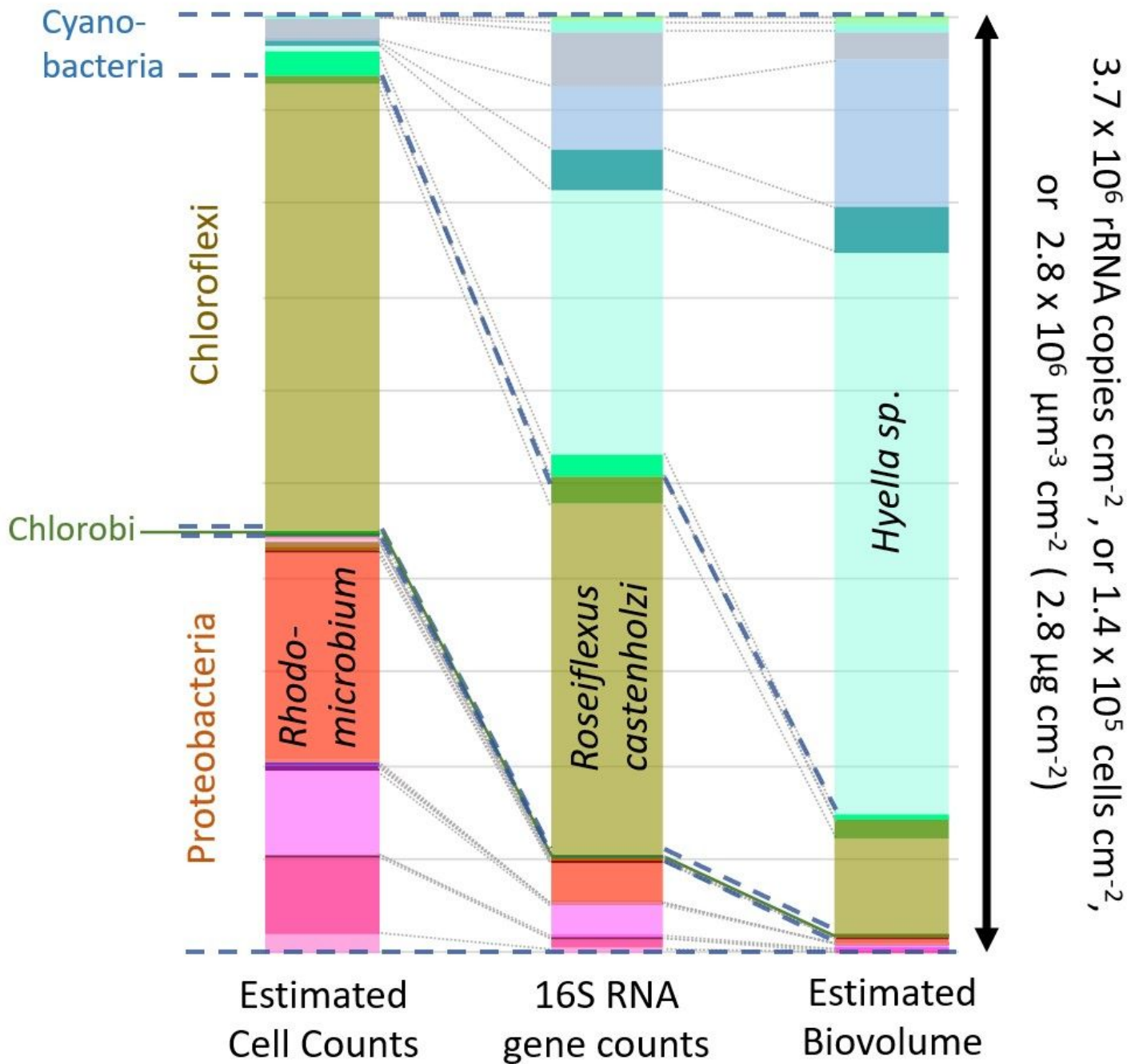


Figure 2

Microbial community structure as determined by ribosomal counts, and as estimated on the basis of cell numbers and biovolume. Legend: Estimation of microbial community structure based on experimental ribosomal counts (central column), estimated cell number (left column) and estimated biovolume (right column) in a single, exemplary dataset using allometric corrections based on Eq. 1. The dataset is from Roush et al. (2020) (22), and includes the subset of taxonomically assignable phototrophic bacteria from an endolithic microbiome on coastal marine carbonate rocks. Only three exemplary phototrophs are labeled, but full, taxonomically explicit distributional data are in Table S3 (see

Additional file 3). For ease of comparison, results are graphically presented as relative frequencies, but absolute scales of areal abundance are indicated on the arrow to the right.

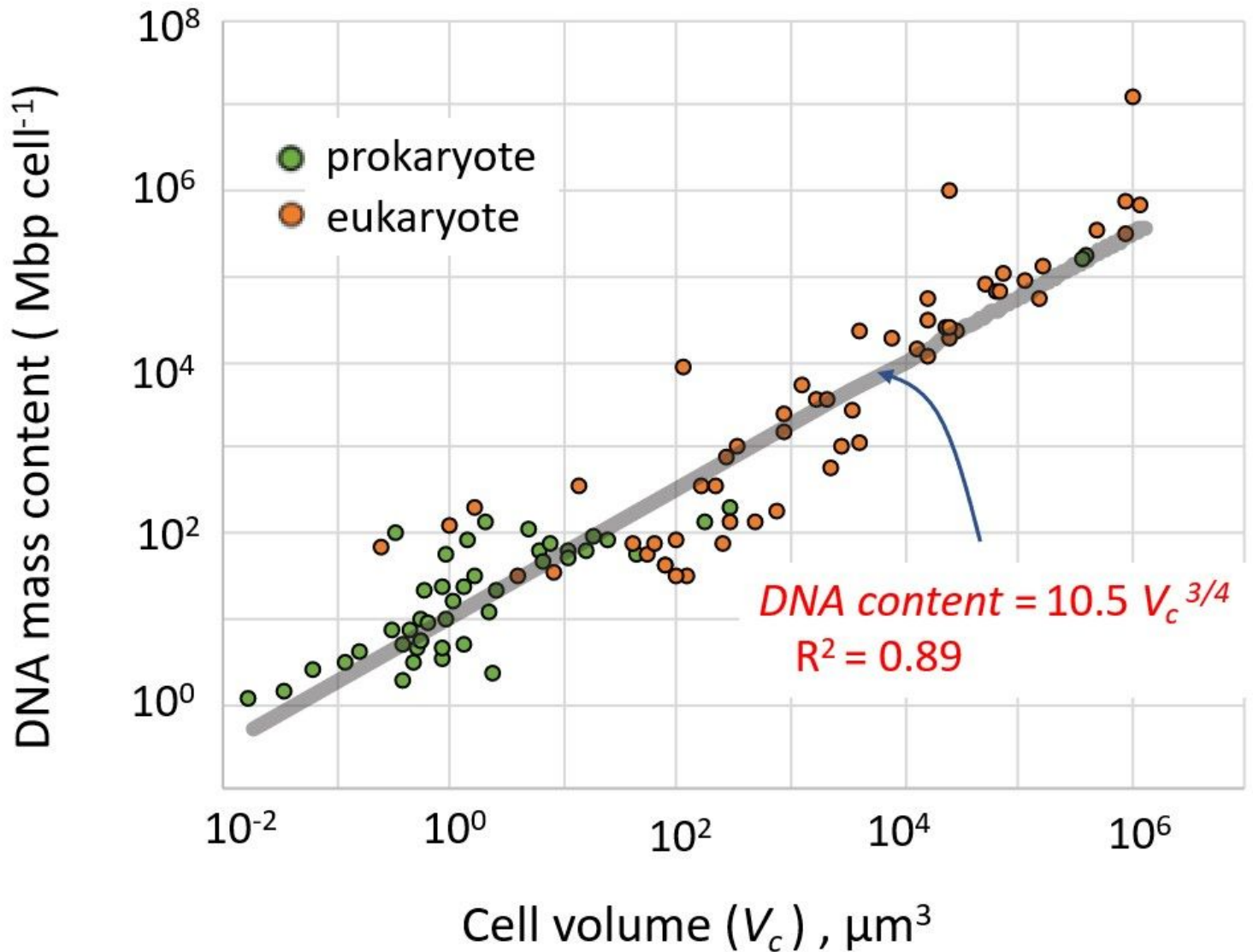


Figure 3

DNA content scales among microbes with cell volume as a power function with an exponent of $\frac{3}{4}$. Legend: Entries are from a subset of those in Table S1 (n = 60, see Additional file 1), and determinations by Shuter et al (35) (n=39). Orange points belong to eukaryotic microbes and green points to prokaryotes. Full statistics for the fit (in red type) are given in Table S2 (Additional file 2).

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Additionalfile1.xlsx](#)
- [Additionalfile2.docx](#)

- [Additionalfile3.xlsx](#)
- [Additionalfile4.docx](#)
- [Additionalfile5.docx](#)
- [Additionalfile6.docx](#)