

Reduced Metagenome Sequencing for strain-resolution taxonomic profiles

Lars Snipen (✉ lars.snipen@nmbu.no)

Norwegian University of Life Sciences <https://orcid.org/0000-0002-2883-861X>

Inga-Leena Angell

Norwegian University of Life Sciences, Faculty of Chemistry, Biotechnology and Food Science

Torbjørn Rognes

University of Oslo, Department of Informatics

Knut Rudi

Norwegian University of Life Sciences, Faculty of Chemistry, Biotechnology and Food Science

Methodology

Keywords: metagenome, strains, ddRADseq

Posted Date: January 22nd, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-44151/v3>

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 March 29th, 2021. See the published version at <https://doi.org/10.1186/s40168-021-01019-8>.

1

2 **Reduced Metagenome Sequencing for strain-resolution**
3 **taxonomic profiles**

4

5 Lars Snipen*

6 Department of Chemistry, Biotechnology and Food Sciences, Norwegian University of Life Sciences

7 P.O. Box 5003, NO-1432, Ås, NORWAY

8 Email: lars.snipen@nmbu.no

9 *Corresponding author

10

11 Inga-Leena Angell

12 Department of Chemistry, Biotechnology and Food Sciences, Norwegian University of Life Sciences

13 P.O. Box 5003, NO-1432, Ås, NORWAY

14 Email: inga.angell@nmbu.no

15

16 Torbjørn Rognes

17 Department of Informatics, University of Oslo

18 P.O. Box 1080 Blindern, NO-0316, Oslo, NORWAY

19 Email: torognes@ifi.uio.no

20

21 Knut Rudi

22 Department of Chemistry, Biotechnology and Food Sciences, Norwegian University of Life Sciences

23 P.O. Box 5003, NO-1432, Ås, NORWAY

24 Email: knut.rudi@nmbu.no

25

26 **Abstract**

27 **Background**

28 Studies of shifts in microbial community composition has many applications. For studies at species or
29 subspecies levels, the 16S amplicon sequencing lacks resolution, and is often replaced by full shotgun
30 sequencing. Due to higher costs, this restricts the number of samples sequenced. As an alternative to a
31 full shotgun sequencing we have investigated the use of Reduced Metagenome Sequencing (RMS) to
32 estimate the composition of a microbial community. This involves the use of double-digested restriction
33 associated DNA sequencing, which means only a smaller fraction of the genomes are sequenced. The
34 read sets obtained by this approach have properties different from both amplicon and shotgun data,
35 and analysis pipelines for both can either not be used at all or do not explore the full potential of RMS
36 data.

37 **Results**

38 We suggest a procedure for analyzing such data, based on fragment clustering and the use of a
39 constrained ordinary least square de-convolution for estimating the relative abundance of all
40 community members. Mock-community data sets shows the potential to clearly separate between
41 strains even when the 16S is 100% identical and genome-wide differences is <0.02 , indicating RMS has a
42 very high resolution. From a simulation study we compare RMS to shotgun sequencing and show that
43 we get improved abundance estimates when the community has many very closely related genomes.
44 From a real data set of infants guts we show that RMS is capable of detecting a strain-diversity gradient
45 for *Escherichia coli* across time.

46 **Conclusion**

47 We find that RMS is a good alternative to either metabarcoding or shotgun sequencing when it comes to
48 resolving microbial communities at the strain-level. Like shotgun metagenomics, it requires a good
49 database of reference genomes, and is well suited for studies of the human gut or other communities
50 where many reference genomes exist. A data analysis pipeline is offered, as an R package at
51 <https://github.com/larssnip/microRMS>.

52

53 **Keyword**

54 Metagenome, strains, ddRADseq

55

56

57 **Background**

58 The study of microbial communities relies on the sequencing of microbial DNA, and current practice can
59 be divided into two main approaches: Metabarcoding, also known as amplicon or targeted sequencing,
60 and shotgun sequencing of random fragments from the entire genome [1]. The amplicon-approach is
61 primarily used for revealing the taxonomic composition, but may also be used to study the distribution
62 of targeted functional genes [2]. Shotgun sequencing provides a potentially more detailed information
63 about the community genomes, the microbiome, and is typically used for studies that digs beyond the
64 composition and into the genomic function. Shotgun microbiome sequencing requires significantly more
65 efforts in sequencing, data processing and analysis compared to meta-barcoding.

66 In most microbiome studies the composition is of interest, and in some cases it is all we require. Shifts in
67 composition may be used as indicators in various ways, e.g. microbiota profiles in forensics [3], or in the

68 surveillance of environments [4]. For communities like the human gut, extensive studies of the
69 composition has given us the big picture, but recent investigations indicates that differences at the
70 strain level may be crucial for phenotypic differences [5, 6]. Common to these problems is the need for
71 high-resolution taxonomic profiles that can be collected with moderate efforts and in a reproducible
72 way. Such studies often require many samples in order to capture the biological variation, and since
73 sequencing and computational resources are always limited, the simpler amplicon approach is often
74 preferred to a deep shotgun sequencing in order to get enough samples covered. However, the standard
75 approach using the 16S rRNA gene marker has a limited resolution. If separation at the species or strain
76 level is required, the 16S marker is in general too conserved, and a full shotgun sequencing seems
77 necessary.

78 As an alternative to do a full shotgun metagenomic sequencing, the use of restriction enzymes to reduce
79 the genomic sequence space has also been employed to investigate microbial communities [7-9]. The
80 double-digested restriction associated DNA sequencing (ddRADseq) idea [10-12] has been along for
81 some time, but its use for metagenome studies is quite new. The main advantage for this approach has
82 been to reduce the sequencing efforts, and thereby costs per sample. In short, this means cutting DNA
83 into fragments using two different restriction enzymes, followed by a PCR amplification and sequencing
84 of the resulting amplicons. This procedure falls between the full shotgun sequencing and the classical
85 use of a specified marker gene (16S). It has some resemblance to shotgun sequencing, since from each
86 genome we sequence many different, in some sense random, fragments that varies in size and number
87 between genomes. However, this also resembles metabarcoding since multiple copies of a certain
88 genome will produce the exact same fragments, and the reads are from these fragments each time.
89 Thus, the approach has been termed Reduced Metagenome Sequencing (RMS). The wet-lab protocols
90 for ddRADseq are well established [8].

91 There are some difficulties that arise with the RMS approach. When target sequencing a pre-defined
92 marker gene like the 16S, we may cluster reads into OTUs or sequence variants, where each cluster
93 represents some taxon. RMS reads may also be clustered, but each taxon give rise to a variable number
94 of distinct fragments, and it is difficult to infer the taxonomic composition from such read clusters
95 without mapping to some references. Also, due to the variable lengths and compositions of the
96 fragments, the PCR-amplification efficiency must be expected to vary and create biases. For purely
97 predictive purposes, such reference-free approaches may still be useful, as suggested by [9].

98 In this paper we focus on the use of RMS data for estimating taxonomic profiles, as an alternative to a
99 full shotgun sequencing. This means reads are in some way mapped to a database of reference
100 genomes, often referred to as closed-reference analysis. Our focus is on the estimation of high-
101 resolution abundances, i.e. species or strain-level profiles, which is not attainable by conventional 16S
102 sequencing. It is possible to use computational tools designed for shotgun data directly in the RMS
103 setting. This may produce helpful results, but do not utilize all the information we have in this case. We
104 propose an alternative analysis approach, using fragment clustering and a constrained least squares
105 estimation. Based on mock community data and simulations, we demonstrate some important aspects
106 of RMS data, and show the potential for RMS to improve composition estimates at the strain level. We
107 also include an example of using RMS to estimate strain-diversity for *Escherichia coli* in the infant gut
108 microbiome. The analysis tools, along with some tutorial material, is freely available as an R package at
109 the GitHub site <https://github.com/larssnip/microRMS>.

110

111 **Results**

112 We have explored the Reduced Metagenome Sequencing approach for studying the composition of
113 microbial communities, with a focus on high-resolution profiles. The RMS idea is to cut genomes into

114 fragments using restriction enzymes, then amplify and sequence the resulting fragments. In this study
115 we have focused on the restriction enzymes EcoRI and MseI. The data processing pipeline illustrated in
116 Figure 1 will apply to any choice of enzymes, but some of the choices made along the way may change.
117 In Figure 2 we see how RMS fragment number and lengths distribute for a selection of species typically
118 found in the human gut. This will clearly change if other restriction enzymes are used. For each species
119 we randomly selected 10 sequenced strains, and all results are based on retrieving the RMS fragments *in*
120 *silico* from the genomes, using the cutting motifs GAATTC (EcoRI) and TTAA (MseI). In the upper panel
121 we notice that the number of RMS fragments per mega base pairs varies a lot between species, but less
122 within each species. The number of RMS fragments is typically limited by the occurrence of the longer
123 cutting motif, in this case GAATTC. Since this is a GC-poor motif, there is an effect of GC-content, and
124 the grey sector indicates where the numbers should have been had it been random DNA. In the lower
125 panel the densities show how fragment lengths distribute. Here we selected three species only, having
126 low (*F. nucleatum*), medium (*E.coli*) and high (*B. longum*) GC-content. The fragment lengths are typically
127 governed by the occurrence of the shorter motif, in this case TTAA, and again there is a huge effect of
128 GC-content. The GC-poor *F. nucleatum* has very short fragments, since the short motif occurs more
129 frequently than in the more GC-rich species. There are fragments longer than 1000 bases, but these
130 typically produce low signals after PCR amplification and are not shown here. For any choice of
131 restriction enzymes, similar investigations should be made to see how many and how long fragments
132 one should expect from the species most likely to be found in the targeted samples.

133 Figure 3 illustrates the potential for high-resolution using the RMS method. Here we have considered
134 the 27 genomes used in our mock community below. In panel A we used the 16S sequence from each
135 genome. These were aligned (MUSCLE, [13]) and the p-distance (1 minus identity, i.e. distance 0.01
136 corresponds to 99% identity) between them computed using the ape-package in R [14]. We notice that

137 strains within the same species are identical or close to identical, and even the two *Staphylococcus*
138 species are difficult to separate, with p-distance <0.01. OTUs based on 16S data are usually clustered at
139 distance 0.03. In panel B we computed the p-distances based on whole genomes, using the fastANI
140 software [15]. This separates the *Staphylococci* better than 16S, but strains within the same species are
141 again quite similar. The two strains of *Helicobacter pylori* have p-distance 0.04 between them, indicating
142 96% of their genomes are identical. In panel C is the correlation distance between genomes based on
143 RMS fragment copy numbers. From the genomes we get the copy number matrix X , with one row for
144 each fragment cluster and one column for each genome. The two *L. gasseri* strains have a small
145 correlation distance, making them as good as impossible to separate with RMS. But, the other species
146 with multiple strains are far better off, and the correlation distance of 0.22 between the two *S. mutans*
147 strains should be large enough for discrimination between them.

148 Previous use of ddRADseq have reported biases in the signals, most likely introduced in the PCR
149 amplification of such fragments [16]. To investigate this we used mock data from [8], where some
150 samples included a single genome only. In Figure 4 we have plotted the data from four such samples.
151 From panel A we clearly see how fragment length affects the relative read count signals by the banana-
152 shaped cloud of strong signals. The cloud of very low signals (note log-transformed y-axis) is due to
153 noise. In panel B there seems to be no bias due to GC-content. Based on repeated observations of
154 similar patterns, we propose a length-normalization of the RMS signals. In panel C we show the effect of
155 this procedure described in the Methods section. We also suggest, based on panels A and C, that only
156 fragments within the length interval 30-500 bases should be used in the downstream analyses,
157 highlighted in panel C of Figure 4. In this interval the length bias is small, and the normalization
158 procedure will not affect the data very much, which is always a good thing. From Figure 2 we also saw
159 that most fragments are in this length-interval when using the current restriction enzymes. Clearly,
160 these limits must be re-considered if other enzymes are used.

161 Next, we used the RMS approach on the mock community data. In Figure 5 we show a classical stacked
162 bar plot displaying the estimated composition of the 20 genome mock. The estimates are based on the
163 constrained ordinary least squares (COLS, see Methods for details) procedure, with 0.1 trimming,
164 described in the Methods section. The database contains the genomes from all the 27 genomes, but
165 only the original 20 genomes were included in this sample. Proportions are well estimated, and of the
166 seven extra genomes absent in the samples, six of them are correctly estimated to have no
167 contributions. The only false positive (weak) signal is from *L. gasseri* ATCC 33323, as expected from the
168 results in Figure 3. In Figure 6 we display the actual versus the predicted relative abundances as
169 scatterplots for each of the other mock-combinations. Here the extra strains were spiked-in, one by one,
170 and in one sample all seven were added. From this we note that proportions are in most cases well
171 estimated, where strains who are absent are estimated with zero proportions, and when strains are
172 spiked-in they are estimated with fairly accurate proportions. There are two exceptions. As seen also in
173 Figure 5, the *L. gasseri* ATCC 33323 strain is estimated as (weakly) present also in those cases it is
174 absent, and the spiking-in of *H. pylori* NCTC 11637 seems to have failed in some way, coming out with
175 much too low abundance in the two cases where it is present.

176 We also tested the RMS approach against a standard shotgun sequencing procedure using simulated
177 data. We focused on human gut-like communities of three different resolutions, where community
178 members have a minimum whole-genome p-distance of 0.05, 0.02 or 0.01. This resulted in 291, 601 and
179 1086 community genomes, respectively. In all cases the samples contained reads from 100 randomly
180 sampled present genomes, but the databases (RMS and Kraken2) contained all community genomes,
181 both those present and absent. Both from RMS and shotgun data we re-estimated the relative
182 abundance of every single genome in the database, using the COLS method for RMS data and Kraken2
183 [17] with a custom database for the shotgun data. To evaluate the results, we computed the Manhattan
184 distance (or L_1 norm) between actual and predicted relative abundances, as suggested in [18]. Thus, a

185 Manhattan distance of $D=0$ means we estimate all relative abundances perfectly. In Figure 7 the actual
186 versus the predicted abundances are plotted as scatterplots. We observe, as expected, that predictions
187 are poorer for lower p-distances, i.e. it becomes more difficult to distinguish genomes as they become
188 more similar. However, the difference between RMS (upper panels) and shotgun (lower panels) data is
189 striking. With the RMS approach we can estimate the abundance of each genome quite well, while for
190 shotgun data the variance becomes huge for the highest resolutions, with predicted abundances up to
191 three times larger or smaller than the actual abundances.

192 Finally, we include some real data to illustrate our use of RMS. First, we use it for quantifying strain-
193 diversity of *E. coli* in the infant gut. At the time of writing, there are 1066 complete *E. coli* genomes in
194 the RefSeq database (<https://ftp.ncbi.nlm.nih.gov/genomes/refseq/>). This is an example of a genome
195 collection where we find some very similar genomes. We computed the whole-genome p-distance
196 between all pairs of genomes, using the MASH software [19], as well as the RMS correlation distances
197 described in the Methods section. In Figure 8 the left panel A scatterplot indicates how these distances
198 relate to each other. We only plot the distance to the nearest neighbor for each genome, and the grey
199 dots are for all 1066 genomes. Note that some of distances are zero (both distance measures), indicating
200 RefSeq contains multiple copies of identical genomes. The scatterplot also relates the correlation
201 distance to the more familiar p-distance, and we observe that a correlation distance around 0.30 here
202 corresponds to a p-distance of roughly 0.01, i.e. genomes of 99% identity.

203 Using all 1066 genomes turns out impossible, because the copy number matrix produces an infinite
204 condition value. This is due to the identical genomes, which is (theoretically) impossible to separate.
205 Hence, we employed the genome clustering described in the Methods section. Setting the maximum
206 tolerated condition value at 100 produced 54 genome clusters, i.e. the *E. coli* population is divided into
207 54 subgroups with this resolution. The black dots in Figure 7 are the nearest neighbor distances for

208 these 54 genomes. In the right panel B we indicate where these are located in a neighbor joining tree
209 based on the p-distances between all strains.

210 Six of the samples from the infant guts were also subject to a conventional shotgun sequencing. We first
211 made a comparison between shotgun and RMS, using Kraken2 and a custom database for assigning
212 reads to the exact same 54 *E. coli* genomes that we used for the RMS analysis. Since we do not know the
213 true composition of these samples, we only considered which strains were estimated to be present or
214 absent in a sample. In the left panel A of Figure 9, we show a tree of the 54 strains, based on whole-
215 genome p-distances, where we have colored the leaf nodes by how they were classified by either
216 shotgun or RMS data in one of the samples. We can see that from the shotgun data, and the Kraken2
217 assignments, 51 of the 54 strains were assigned reads, and thereby being present (grey or black), while
218 the RMS results only estimate 16 strains as present (black only), of which half is from the same clade at
219 the top of the tree. The other five samples show a similar trend: The shotgun approach will assign reads
220 to a majority of strains, while the RMS approach is more specific, stating fewer strains are present.

221 In the right panel B of Figure 9 the boxplot shows the number of strains found present by RMS in all 94
222 infants at 4 different times after birth. The variance is large, partly due to biological variation. Still, the
223 trend of a growing diversity by age is clear, and a simple ANOVA analysis confirmed a highly significant
224 increase in diversity from birth (Meconium) to all later times, especially to 12 months.

225 As a second illustration of the use of RMS and our deconvolution method, we re-analyzed the data from
226 [7]. In their paper they sequenced 3 human gut samples with both a conventional shotgun approach as
227 well as an RMS approach. They used other restriction enzymes than we did in our analyses above, the
228 *NlaIII* with cutting motif `CATG` and *HpyCH4IV* with `ACGT`. In [7] all data were profiled by MetaPhlan2
229 [20], i.e. both shotgun and RMS data were treated the same way. In our re-analysis we used Kraken2
230 instead of MetaPhlan2, and in addition we also used our RMS-specific method on the RMS data. As a

231 genome database, both for Kraken2 and our own method, we used the MGnify collection of human gut
232 genomes [21]. This consists of 4644 genomes isolated from human gut samples, and clustered at 95%
233 identity, i.e. each genome represents a cluster of genomes with more than 95% identity.

234 In the left panel A of Figure 10 we show a principal component analysis (PCA) plot of the three samples
235 based on shotgun sequencing and Kraken2 profiling (shotgun+kraken), RMS sequencing and Kraken2
236 (ddRADseq+kraken2) and RMS sequencing profiled by our COLS method (ddRADseq+COLS). Only the 100
237 overall most abundant taxa were included to make the figure, and their relative abundances were
238 transformed by the centered log-ratio transform [22] prior to the PCA computations and plotting.

239 In [7] the most abundant species identified by MetaPhlan2 in these samples was *Faecalibacterium*
240 *prausnitzii*. This is a common human gut species, but is also known to consist of multiple phlotypes,
241 where different phlotypes have been reported to be associated with differing disease developments
242 [23, 24]. Thus, it is of some interest to separate between variants of this species. In the MGnify
243 collection we find 9 different genome clusters named *F. prausnitzii*. We therefore used our
244 deconvolution approach to separate between them, as plotted in the right panel B of Figure 10.

245

246 **Discussion**

247 The upper panel of Figure 2 shows that, using the restriction enzymes of this study, the RMS fragment
248 density varies by GC-content. Most genomes also have fewer fragments than expected in random DNA,
249 indicating a negative selection of the cut sites. However, multiplying by genome size, we find that most
250 genomes have in the range of hundred to thousand fragments. The lower panel of Figure 1 shows most
251 fragments are rather short. This is a good thing, since longer fragments amplify poorly, and we found
252 that by only focusing on fragments in the length-interval 30-500 bases, we obtain strong signals without

253 too much PCR bias. These results apply to our chosen restriction enzymes and should always be
254 investigated for any alternative choices of enzymes.

255 From Figure 3 we clearly see the potential for RMS to resolve strains at a level which is impossible with
256 16S, and even difficult with shotgun sequencing. In Panel A we used full-length 16S sequences, but still
257 the separation is very poor between closely related genomes. The distances in panel B reflects how
258 similar the genomes are in overall nucleotide identity. As expected, strains within a species have p-
259 distance less than 0.05, i.e. more than 95% identical. Panel C demonstrates that even closely related
260 strains have rather large correlation distance, indicating a good number of unique RMS fragments. This
261 may seem strange, how can genomes be so similar in p-distance, but still have different RMS fragments?
262 Mutations in restriction cut sites as well as re-arrangements of genomic regions will both create/destroy
263 RMS fragments, but have little impact on the whole-genome distance. Only if two or more genomes
264 share the vast majority of RMS fragments, we would see a small correlation-distances, and a shallow
265 branch in the dendrogram. From panel C in Figure 2 we expect to be able to separate all genomes by
266 RMS, except perhaps the two *L. gasseri* strains.

267 Since the RMS approach involves a PCR step, we must expect some biases. In Figure 4 (panel A) we
268 observe a distinct effect of fragment length on the relative signals strengths we get, based on single-
269 genome samples from a previous study. However, the GC-content of the fragments does not seem to
270 have any effect (panel B), unlike what was reported by [16]. The lower panels of Figure 4 illustrate our
271 proposed way of handling the length-bias. First, only use fragments in the length interval 30-500 bases,
272 highlighted in brown color in panel C. As we saw in Figure 2, short fragments account for the vast
273 majority anyway, and we have found that 80-90% of the reads will map to fragments in this interval.
274 Next, we propose a simple normalization, illustrated in panel C. The cloud of strong signals is
275 straightened. Note that due to the log-transformed y-axis it looks like weak signals (noise) are heavily
276 distorted by the normalization, while their values actually change very little. If other restriction enzymes

277 are used, the fragment lengths may be different, and the limits of 30-500 should be reconsidered.

278 However, the length bias corrected by the normalization will probably be of the same type, since this is a

279 PCR effect which is independent of the restriction enzymes.

280 The mock data results in Figure 5 and 6 reveal that with the RMS approach and the COLS algorithm we

281 can estimate relative abundances fairly well. It should be noted that the actual abundances are probably

282 not exact, as they rarely are in experimental data. The mock composition was designed by 16S copies,

283 and the transformation to genome copies is not without uncertainty, since most of these species are

284 known to have variable 16S copy numbers. Most important is that strains who are absent from a sample

285 are also estimated to zero abundance, i.e. they show up in the lower left corner of the panels in Figure

286 6. The exception is *L. gasseri* ATCC 33323. This is simply too similar to the other *L. gasseri* strain, as seen

287 in Figure 3 as well. Only five fragment clusters are unique to the ATCC 33323 strain, and with some noise

288 signals on some of these, it appears to be present even when it is not.

289 In virtually all cases the three replicates (marker types) of each sample show very similar results,

290 indicating there is very little variance in the RMS procedure as such. Hence, any deviations between

291 actual and estimated proportions are most likely due to some systematic effect. On closer inspection we

292 found that the RMS data display what we denote as a fragment bias. Fragments unique to a genome

293 should in theory all produce similar read counts, with some variation due to the randomness of

294 sequencing. This is not the case. Some fragments consistently produce strong signals and others weak.

295 This is also remarkably stable across all samples where a particular genome is present. We accounted for

296 this in our simulation study, adding a random scaling to all fragments, and supplementary Figure 2

297 shows the distribution of this fragment bias. So far we have failed to reveal the cause of this effect. If we

298 could understand and compensate for it, this would improve the precision and thereby the resolution of

299 the method even further.

300 We used simulated data to compare the RMS approach to the use of shotgun sequencing in combination
301 with the Kraken2 tool for re-estimating the relative abundances of each genome. Kraken2 is only one
302 out of several tools for estimating metagenome composition, but we chose this because it has a good
303 reputation, will always try to classify reads to the genomes of its database, but most importantly, the
304 genomes in the database can be easily customized. To make the comparison fair, the database must be
305 identical for both the RMS and shotgun approach. Using a generic database is bound to produce poorer
306 results compared to one where the exact genomes under study are in the database. An alternative tool
307 like MetaPhlAn2 [20] assigns reads no lower than to the species level, but has the extension StrainPhlAn
308 [25] for a strain level analysis. It is, however, difficult to compare StrainPhlAn output to the ones we get
309 here, since we focus on relative abundances of *a priori* defined genomes, while StrainPhlAn identifies
310 strains *a posteriori* by aligning reads to a set of marker genes and output a multiple sequence alignment.
311 It seems to us these are two quite different approaches to a strain-level analysis.

312 In Figure 7 we show some results of our simulation study. The left panels are from a community where
313 no members have a p-distance below 0.05 to another member. This is roughly a community with one
314 genome from each species. Here both methods perform extremely well, and the shotgun+Kraken2 is the
315 best, with Manhattan distance $D=0.034$. However, as the communities (and databases) are filled up with
316 more and more similar genomes, the picture changes (middle and right panels). The shotgun+Kraken2
317 results are getting dramatically poorer, with highly fluctuating estimates of relative abundances. The
318 RMS+COLS approach is also poorer, but not nearly as bad. While the Kraken2 results seem to be fairly
319 unbiased, but with a huge variance, the COLS results have a small variance, at the cost of some bias in
320 giving weak abundance to some absent genomes, seen in the lower left corners of the 98% and 99%
321 panels. Our explanation for these results is that with shotgun sequencing most reads will match multiple
322 genomes in the database, and Kraken2 will then assign to the lowest common ancestor, i.e. the species.
323 Thus, species abundances become extremely precise, but too few reads are left at the strain level to get

324 reliable estimates. In our COLS algorithm for RMS data, we also have many fragment clusters who are
325 present in multiple genomes, but since we have the copy number matrix with this exact information, the
326 constrained least square solution spreads the signal across all genomes instead of assigning it to their
327 common ancestor. It should be mentioned that this idea has some resemblance to what was proposed
328 by [26], using methods from RNAseq data as an alternative approach for analyzing shotgun data. The
329 Kraken2 software also has an extension in the Bracken software [27], re-estimating low-rank
330 abundances based on the higher-rank assignments, but is difficult to use below the species rank.

331 There is, as for any method, a limit to the resolution obtained by RMS. The example with 1066 *E. coli*
332 genomes illustrates this. Many of these are more than 99% identical, some even 100%, as seen in Figure
333 8. We plotted the correlation distance between all genomes against the p-distance for the same pairs, to
334 illustrate how they are related. An RMS correlation distance of around 0.30 corresponds roughly to a p-
335 distance of 0.01 (99% identity) in this case. We employ a genome clustering, where we only keep a
336 selection of the genomes, ensuring a minimum difference between them. This means each cluster
337 centroid represents a subgroup of highly similar strains. When using the COLS algorithm the resolution is
338 limited by the condition value of the fragment copy number matrix. A very large condition value
339 indicates the estimated abundances will be unstable. Condition values of 10^2 , 10^3 or even 10^4 may be
340 used to obtain a gradually higher resolutions, but at the cost of more uncertain results. Even with the
341 lowest threshold at 10^2 we get 54 subgroups in the analysis, and it is likely that in many cases such a
342 resolution will suffice. As seen in Figure 8, these strains are typically 98-99% identical (p-distance 0.02-
343 0.01) and represent the full tree of all strains quite well.

344 A shotgun sequencing should in theory be able to separate anything below 100% identical, but in
345 practice not. Reads are not without errors, and read coverage is often poor for low-abundance taxa. The
346 results in Figure 9 underlines this. The shotgun data indicate almost all genomes in the database are
347 present in the sample. With the RMS approach, much fewer genomes are present. In one out of the six

348 comparable samples the methods came out with the exact same genomes as present and absent. In the
349 others the shotgun data always results in more detected strains. It is reasonable to suspect both
350 methods are too sensitive, assigning too many subgroups as present, but RMS seems far better in this
351 respect. The total fraction of *E. coli* is small in these samples (around 1%) but the absolute number of
352 reads assigned to this species are in the same range for both methods (around 1000). It is in fact slightly
353 larger for the RMS data, hence the increased prevalence from shotgun data is not due to increased
354 coverage. For shotgun data reads can originate from all locations on the genomes, making it notoriously
355 difficult to map a read correctly when genomes are as similar as here, and given that reads may contain
356 errors. RMS reads are assigned to the *a priori* known fragments and allows for some slack due to
357 sequencing error. Also, if genome A and B share 50% of their fragments, but only the genome A
358 fragments have signal, the COLS algorithm will assign abundance 0.0 to genome B even if 50% of its
359 fragments have signal. This is possible because we know these fragments are shared with genome A,
360 and since the unique genome A fragments have signal while the unique genome B have not, the shared
361 fragment signals are all allocated to genome A, giving no abundance to genome B.

362 The boxplots in panel B of Figure 9 is an example of how we use RMS to detect a change in strain-
363 diversity over time in the infant gut. The increasing diversity by age is as expected. This example also
364 illustrates how patterns emerge because we were able to sequence many samples, rather than deep
365 sequencing of a few, where the biological variation probably would obscure the results. Such a high-
366 resolution analysis would not be possible by 16S analysis.

367 The re-analysis of the data from [7] is an example of using completely different restriction enzymes. The
368 two four-base cutters result in far more fragment per genome than we saw in Figure 2, but apart from
369 this the analysis we did were identical to what we have done above. The left panel of Figure 10 shows
370 that shotgun data and RMS data, assigned both by Kraken2 and our algorithm, results in the same big
371 picture. The difference between methods is small compared to the difference between samples, which

372 is the same conclusion reached by the authors of the original paper. In the right panel, however, we
373 show that with our RMS-specific deconvolution we may now also estimate abundances below the
374 species level. The original results, using MetaPhlan2, did not dig beyond the species level, but this
375 seems to some degree possible with the approach we have suggested in this paper. The species *F.*
376 *prausnitzii* is exactly a species where such analyses may be of some interest. We observe some
377 difference in strain abundances here, but three samples is of course far too few to reach any conclusion
378 along this road.

379 The RMS has been proposed as a low-cost alternative to a full shotgun sequencing [7], since we only
380 sequence the amplified fragments accounting for a fraction of the entire genomes. This is true if you use
381 a reference-free approach where you need to cluster the reads, and hence need to have sequenced the
382 same region of a genome several times in order to say something about abundance. However, as long as
383 reads are mapped to reference genomes, this difference in library complexity is less important. Instead,
384 the potential gain in using RMS lies in precise estimates of strain-resolution profiles. As for shotgun data,
385 there is no theoretical lower sequencing depth that is required, the more reads the better. For the mock
386 data results in Figures 5 and 6, where strains separated nicely, each sample had between 1 and 2 million
387 reads mapped to some fragments, resulting in mostly 10-100 read per fragment. This we consider a very
388 good coverage. As always, high coverage is needed for detecting low abundance taxa, but is not in itself
389 required for separating closely related strains. A bottleneck for RMS is the fragment bias previously
390 mentioned. For some reason, fragments from the same genomes tend to get quite different read
391 counts, in a reproducible way. If a genome has as very few fragments, the average read count for these
392 is not as stable as with many fragments.

393 We believe our results indicate the RMS approach for metagenome profiling is something to explore
394 further. We have focused a lot on one pair of restriction enzymes in this study, but other enzymes are
395 used for similar studies [7, 9]. The choice of enzymes will affect the number and length of fragments, but

396 apart from this the data analysis procedure we propose here may be used, as we illustrate by the re-
397 analysis of the data from [7]. In the supplied software (R package) there are options for using any pair of
398 restriction enzymes. The RMS approach, like the shotgun metagenome approach, requires sequenced
399 reference genomes to map against in order to produce taxonomic profiles. To obtain this at the strain
400 level, we need good reference databases. The good news is that recent extensive efforts provide us with
401 many new reference genomes, especially for the human gut[28-32]. We believe that with evolving
402 sequencing technologies, the quality of metagenome assembled genomes (MAGs) will improve
403 drastically, and the road lies open for more strain-level profiling.

404

405 **Conclusion**

406 We have demonstrated that the RMS approach can be used for profiling of microbial communities down
407 to the strain level. Compared to the conventional 16S approach, we find that strains with identical 16S
408 genes are clearly discriminated by RMS, and we can estimate abundances for such strains in the same
409 sample. The reason for this is simply that even genomes with identical 16S sequences will in most cases
410 differ in a fair number of RMS fragments, enough to obtain strain-specific signals for the COLS algorithm.
411 Compared to the shotgun metagenome approach, the RMS offers an advantage in only sequencing *a*
412 *priori* known amplicons, and we may construct a copy number matrix revealing the relations between all
413 reference genomes prior to any sequencing. From this information, and the suggested constrained
414 ordinary least squares estimation algorithm, we can obtain strain-level abundance estimates at least as
415 good as the popular metagenome tool Kraken2. A clustering of genomes into species subgroups is
416 proposed, as a way of balancing high resolution against precision in estimated abundances.

417 Based on this, we conclude that the RMS approach is worth pursuing, as a tool for studies of
418 composition in the human gut or other microbial communities of particular interest and where a
419 comprehensive collection of reference genomes exists. An R-package with the data analysis methods
420 suggested here, as well as tutorials, is available in GitHub at <https://github.com/larssnip/microRMS>.

421

422 **Methods**

423 **Mock data**

424 In order to test the RMS approach, and learn about how such data behave, a mock-community study
425 was conducted. As a basis we used a mock community of 20 genomes obtained through BEI Resources,
426 NIAID, NIH as part of the Human Microbiome Project (Genomic DNA from Microbial Mock Community B
427 (Even, Low Concentration), v5.1L, for 16S rRNA Gene Sequencing, HM-782D, [33]), see Table 1. This
428 mock has been constructed to yield 100 000 16S copies from each included organism. We converted this
429 into the number of genome copies by dividing 100 000 by the 16S copy number for each organism, as
430 listed in the Ribosomal RNA Database [34]. In addition to this mock itself, we spiked-in 7 additional
431 DSMZ strains (Leibniz Institute DSMZ-German Collection of Microorganisms and Cell Cultures,
432 <https://www.dsmz.de/>). These strains were selected to be highly similar, and with identical 16S gene, to
433 one of the existing strains in the mock, to see if we could separate signals from such closely related
434 organisms. One strain was spiked-in at a time, producing 7 additional samples. The spiked-in genomes
435 were also at 100 000 16S copies, controlled by a droplet digital PCR procedure. Finally, a sample with all
436 27 strains was also used. All these 9 mock mixtures were done in triplicates, resulting in 27 samples. All
437 samples were subject to the wet-lab procedures described in [8] to obtain paired-end Illumina HiSeq
438 reads. The restriction enzymes EcoRI and MseI were used throughout this study. All the strains involved

439 in all samples have whole genome sequence data publicly available, and these were downloaded from
440 the NCBI Genome database (<https://www.ncbi.nlm.nih.gov/genome>).

441

442 Infant gut data

443 As an illustration of a high-resolution analysis, we used a set of RMS data from the gut of infants. The
444 microbiome was sampled from feces of 94 infants at meconium (newborn) and 3, 6 and 12 months age.
445 We used a genome collection consisting of all complete RefSeq genomes of *E. coli* (1066 genomes) in
446 order to look at strain diversity in these samples. Six of the samples were also sequenced by
447 conventional shotgun sequencing, for comparison. All RMS samples were subject to the wet-lab
448 procedures described in [8]. Both RMS and shotguns samples were sequenced by Illumina HiSeq,
449 resulting in 150 bp paired end reads.

450

451 Fragment copy number matrix

452 There exists a number of computational tools for estimating the taxonomic composition of a community
453 based on shotgun data, e.g. Kraken2, MetaPhlan2, CLARK, Kaiju [17, 35-37]. Common to all is that reads
454 are somehow mapped to some database of reference genomes. This is also required for RMS data.
455 Given the reference genomes, and the cutting patterns of the restriction enzymes used (EcoRI and
456 MseI), all RMS fragments were collected *in silico* from each genome. The RMS fragments are simply all
457 sub-sequences starting with an EcoRI motif GAATTC (5' of fragment), and ending by the first
458 downstream MseI motif TTAA, containing none of these motifs inside. Genomes will in general have
459 many such RMS fragments of highly variable lengths.

460 Fragments from closely related genomes may be identical or very similar. Also, some fragments may
461 occur multiple times within a single genome. For this reason we clustered the fragment sequences into
462 fragment clusters, using the VSEARCH software [38] and some specified identity threshold, similar to
463 OTU-clustering for 16S data. An identity threshold of 0.99 was used in this study, but other thresholds
464 were tested without significant changes in results. Each fragment cluster was represented by its centroid
465 sequence and the fragment cluster copy number was stored in a *copy number matrix*. This matrix X has
466 one row for each fragment cluster, and one column for each genome, and the integer in cell $X(i,j)$ is the
467 number of fragments from genome j that belongs to cluster i . If we have a large collection of genomes,
468 this matrix becomes huge. However, most fragment clusters occur in only one, or a few, genomes and
469 most cells in the matrix are zero. Thus, the copy number matrix was stored as a sparse matrix data type,
470 allowing most matrix operations but using comparatively little memory. This copy number matrix is an
471 essential ingredient in the estimation of community abundances, as described below.

472

473 Read processing

474 We used the software VSEARCH [38] for all processing of reads. All reads were subject to a quality
475 filtering, keeping only reads with an expected error rate below 0.02. Read pairs were then merged. Since
476 RMS fragments vary in length, some longer fragments produce non-overlapping reads. Thus, non-
477 merged reads were included as single reads, where the R2-reads were reverse-complemented. To
478 maintain the correct per-fragment read count, all merged reads were given a count of 2, while the single
479 reads counts as 1. All reads were then de-replicated to obtain fasta-files of unique reads for all samples.
480 Proper use of the `--sizein` and `--sizeout` options in VSEARCH allows us to work with the smaller
481 set of unique reads without losing any information about actual read abundances.

482 Next, the processed reads from each sample were mapped to the fragment cluster centroids, using
483 VSEARCH and the identity threshold from the fragment clustering (0.99, see above). This produced a
484 read count matrix \mathbf{Y} , with one row for each fragment cluster and one column for each sample.

485

486 Length-normalization

487 We realized the need for correcting the read count signals due to fragment length PCR bias. First, let \mathbf{y}_k
488 denote raw read counts from sample k , i.e. column k in \mathbf{Y} . Thus, $\mathbf{y}_k(i)$ is the raw read count for fragment
489 cluster i , and L_i is the length of cluster centroid i . Setting aside all clusters with zero signal, the
490 $c_i = \log_{10}(\mathbf{y}_k(i))$ is simply the log read count. Next, we fitted a locally weighted scatterplot smoother (loess)
491 $S(c_i/L_i)$ to these data, thus $S(c_i/L_i)$ is a smooth curve describing how log-read counts c_i vary by fragment
492 length L_i . Then, a correction factor for fragment cluster i is given as

$$493 f_i = 10^{(S_{max} - S(c_i/L_i))}$$

494 where S_{max} is the maximum value on the loess-curve. The normalized read count for any fragment
495 cluster is then

$$496 \mathbf{y}_{k(i)}^* = \mathbf{y}_{k(i)} f_i$$

497 This multiplicative adjustment means fragments with zero signal remain zero also after normalization.

498 This normalization is done for each sample separately. If the database contains a huge number of
499 fragment clusters (many genomes), only a random sub-sample of them may be used to fit the loess
500 model, in order to save time and memory.

501

502 Constrained ordinary least squares (COLS)

503 If all fragment clusters were unique to a single genome, the abundance of each genome would naturally
504 be estimated by averaging the read counts for their corresponding fragment clusters. However, many
505 RMS fragment clusters may be found in several genomes, and more closely related genomes will share
506 more fragment clusters.

507 Prior to sequencing we constructed the copy number matrix from the G genomes in the database. This
508 results in C fragment clusters, thus the copy number matrix \mathbf{X} has C rows and G columns. Let $\mathbf{b}=(b_1,$
509 $b_2,\dots,b_G)$ be the proportion of the various genomes in sample k , i.e. $b_j \geq 0$ and $\sum_{j=1}^G b_j = 1$. Then, it is
510 reasonable to assume that

$$511 \quad E(\mathbf{y}_k) = a \mathbf{X}^t \mathbf{b}$$

512 where \mathbf{y}_k are the (normalized) readcounts for each of the C RMS fragment clusters given the data from
513 sample k , \mathbf{X} is the copy number matrix and a is some positive scaling factor relevant for sample k . Thus,
514 the expected signal for fragment cluster i , $E(\mathbf{y}_k(i))$, is proportional to the linear combination of fragment
515 cluster copy numbers and genome abundances.

516 Given this model, the scaled proportions can be estimated using the constrained ordinary least squares
517 (COLS) approach: Find $\mathbf{s}=\mathbf{a}\mathbf{b}$ that minimizes

$$518 \quad f(\mathbf{s})=(\mathbf{y}_k-\mathbf{X}^t\mathbf{s})^t (\mathbf{y}_k-\mathbf{X}^t\mathbf{s}) \quad \text{where } s_j \geq 0 \quad (1)$$

519 From the requirement that the \mathbf{b} 's must sum to 1.0 we get that $\sum_j s_j = a$ and the estimated relative
520 abundance of each genome is

$$521 \quad \hat{\mathbf{b}} = \frac{\hat{\mathbf{s}}}{\sum_j s_j}$$

522 In the implementation of this de-convolution, we have added the possibility of a trimmed estimate. This
523 means a two-stage estimation procedure: After the initial fitting of the model, as described above, the

524 residuals $\mathbf{y}_k - \mathbf{X}'\hat{\mathbf{s}}$ for all fragment clusters are computed. Then, a user-selected fraction of the fragments
525 with most extreme residuals are discarded, and the model is re-fitted on the trimmed fragment set. This
526 makes the estimated abundances less sensitive to extreme signals from some fragments, but also
527 reduces the size of the data set.

528

529 Correlation distance and genome clustering

530 The COLS algorithm also indirectly suggests the maximum resolution possible to de-convolve. If two
531 genomes are very similar, they will share most RMS fragment clusters, and their respective columns in
532 the copy number matrix \mathbf{X} become similar. The *correlation distance* between two genomes is simply 1
533 minus the correlation between their respective columns in \mathbf{X} . Thus, a correlation distance of 0.0 means
534 the two columns are identical, and the genomes share all fragment clusters. A correlation distance can
535 be maximum 2.0, meaning all fragments present in one genome is absent in the other, and vice versa.

536 When solving eq. (1) we need to invert the matrix $\mathbf{X}'\mathbf{X}$, and if two or more columns are too similar this
537 matrix inversion becomes highly unstable resulting in poor abundance estimates. This instability is often
538 quantified by the condition value of $\mathbf{X}'\mathbf{X}$. A perfect condition value of 1.0 means all columns in \mathbf{X} are
539 orthogonal, i.e. no shared fragments. As columns become more and more correlated, the condition
540 value increases. By computing the condition value from \mathbf{X} we get an idea of how solvable this is, prior to
541 any experimental efforts.

542 Instead of trying to estimate the abundance of all genomes, we cluster them into groups, and replace
543 them by the group centroid genomes, as a representative of each group. The centroid is the one with
544 smallest sum of distances to all the others in the same group. This basically means we get fewer columns
545 and rows in \mathbf{X} . We employed a clustering procedure as follows:

- 546 1. Compute the correlation distance between all pairs of genomes from the columns of X .
- 547 2. Compute a single linkage hierarchical clustering of the genomes based on this. This results in a
548 dendrogram.
- 549 3. Each height in the dendrogram corresponds to an alternative clustering. Choose the largest
550 dendrogram height resulting in a copy number matrix with condition value below a user-
551 specified tolerance.

552 In this way, the user specifies a tolerated upper condition value (e.g. 100 or 1000), and genomes will be
553 clustered to the finest resolution not violating this. A larger tolerance value leads to a finer resolution,
554 but also more unstable estimates.

555

556 Simulation study

557 We also included a simulation study, where we compared the RMS approach to shotgun metagenome
558 sequencing at various resolutions. Genome similarity was computed as whole-genome p-distance, i.e.
559 1.0 minus the Average Nucleotide Identity (ANI). A whole-genome p-distance of 0.0 means identical
560 genomes and above 0.3 means very different genomes. Strains from the same species usually have p-
561 distance below 0.05. In most real communities, like the human gut, we must expect some closely related
562 strains, having a p-distances well below 0.05.

563 In [32] 1520 genomes from the human gut were isolated and sequenced. The whole-genome p-
564 distances between all pairs of these genomes were computed using the MASH software [19], and then
565 used to form clusters at three different resolutions: p-distance 0.05, 0.02 and 0.01. The cluster centroids
566 were used as community members. The following procedure was applied to all communities, separately:
567 From a community of G genomes, a sample contained reads from 100 randomly selected genomes, i.e.

568 100 of the G genomes are present, the remaining $G-100$ are absent. Their abundances were
569 exponentially distributed such that the largest abundance was 100 times the lowest abundance
570 (dynamic range of 100), see supplementary figure 1. Let f_1, f_2, \dots, f_G be the relative abundance for each of
571 the G genomes in the community, i.e. 100 of them are positive and the rest are zero, and they all sum to
572 1.0. These values form the actual relative abundances that we later tried to estimate.

573 This was repeated 25 times for each community, forming 25 different samples. Note that for each
574 sample new 100 present genomes were randomly selected from the sub-population, thus different
575 genomes were present/absent in each sample.

576 Reads were simulated using the ART software [39], using Illumina HiSeq 2500 error profiles, resulting in
577 paired-end reads of 150 bases. For each sample we simulated 1 million read pairs, either as a shotgun
578 sample or as an RMS amplicon sample.

579

580 *Shotgun data*

581 The ART software requires the user to supply the reference sequences to simulate from as well as the
582 number of read pairs to generate. In shotgun metagenome sequencing, the probability of a read pair to
583 originate from genome g is proportional to the abundance of the genome multiplied by its size. After
584 fragmentation of the genomic DNA, the reads are sampled from this fragment pool, and larger and more
585 abundant genomes will contribute with more fragments. Thus, if z_g is the size of genome g , we form a
586 weight for genome as

$$587 w_g = f_g z_g$$

588 Given that we sequenced a million read pairs, these were spread out among the genomes by random
589 sampling using the probabilities

590 $p_g = w_g / \sum_{j=1}^G w_j$

591 resulting in read-counts r_1, r_2, \dots, r_G for each genome. Note that genomes with zero abundance get zero
592 reads. Finally, read-pairs were simulated from each genome, given these read-counts, and assembled
593 into a pair of fastq files. This was then repeated for each sample, producing new sets of fastq files.

594

595 *RMS data*

596 Instead of random fragmentation of the genomic DNA, the RMS protocol results in amplicons based on
597 the fragments we get from restriction enzyme cutting. For each genome sequence we collected the RMS
598 fragments *in silico*, again using the EcoRI and MseI restriction enzyme cutting motifs. Next, we have
599 observed two main biases in how the RMS fragments from a given genome contributes to the pool of
600 sequenced amplicons:

601 First, there is a length bias, especially very long fragments are poorly amplified. Let l_{gk} be the factor that
602 scales the amplification of fragment k in genome g . This is a function of fragment length only, and in
603 supplementary figure 2 we show the function we used for simulating this.

604 Second, we have also observed that some fragments are consistently more or less represented in the
605 reads from a given genome. We denote this the fragment bias. Let v_{gk} be this fragment bias factor for
606 fragment k in genome g , i.e. it may scale the amplification of fragment k up ($v_{gk} > 1$) or down ($v_{gk} < 1$).

607 These factors were sampled at random from the distribution in supplementary figure 2, once for each
608 genome, and then used forever after. Both this distribution, as well as the length bias function, were
609 estimated from real RMS data, using the restriction enzymes described above.

610 Together, this means that the fragments from genome g get the weights

611 $w_k = f_g l_{gk} v_{gk}$

612 where $k=1,2,\dots,F_g$ and F_g is the number of fragments in genome g . All fragments, together with their
613 weights, were assembled for all abundant genomes, and the read-count for each fragment/amplicon
614 was sampled at random, again using probabilities $p_g = w_g / \sum_{j=1}^G w_j$.

615 Note that for shotgun data the weights are only affected by genome abundance and size, while RMS
616 data is affected by genome abundance, number of fragments, length distribution and fragment-bias
617 distribution for the present genomes.

618

619 *Databases*

620 The databases contained all G genomes of the community, both the 100 present at various levels and
621 the $G-100$ absent. For each community, all RMS fragments were found in all G genomes, and a copy
622 number matrix was constructed using a 0.99 identity threshold, as described above.

623 For the shotgun data we used the Kraken2 software [17] to obtain relative abundance estimates. This
624 tool has shown good results in several benchmarking studies [40-42], but more importantly, is equipped
625 with excellent facilities for building a custom database. In order to make a fair comparison to the RMS
626 approach, the database of reference genomes must be the same as in the RMS case. Thus, custom
627 Kraken2 databases were constructed, containing all G genomes of the respective communities. Also, the
628 taxonomy was extended correspondingly, to have a taxonomy-id for every single genome, making it
629 possible for Kraken2 to list hits to each genome.

630

631 *Analysis*

632 The analysis of the RMS data was carried out as described above, but without any genome clustering,
633 resulting in an estimate of the relative abundance of every genome in the database.

634 For the shotgun data, Kraken2 and its custom database was used to assign reads to the genomes, using
635 the default confidence level of 0.0. Only reads assigned to the genome level were counted, since this is
636 our focus. The read count for a genome was divided by the genome size (base pairs), to produce the
637 genome signal. Finally, these signals were divided by the total sum of signals, to produce relative
638 abundances for all genomes in the communities.

639

640

641 **Ethics approval and consent to participate**

642 Not applicable.

643

644 **Consent for publication**

645 Not applicable.

646

647 **Availability of data and materials**

648 The mock datasets generated and analyzed during the current study are available in the Sequence Read
649 Archive repository, under the accession PRJNA574678, see
650 <https://www.ncbi.nlm.nih.gov/bioproject/PRJNA574678/> .

651 The computational methods described in this paper are available as an R-package. It is currently
652 available at GitHub, together with some tutorials describing the analysis steps, see
653 <https://github.com/larssnip/microRMS> .

654

655 **Competing interests**

656 The authors declare that they have no competing interests.

657

658 **Funding**

659 This project has been financed by the Norwegian University of Life Sciences and the project DigiSal NFR
660 248792.

661

662 Author's contributions

663 LS proposed the data analysis methods, did all R programming and drafted the manuscript. IA did all lab
664 work related to experiments. TR did all C++ programming related to VSEARCH. KR conceived the idea. All
665 were involved in discussing methods and editing the manuscript.

666

667 Acknowledgements

668 Not applicable.

669

670 References

- 671 1. Breitwieser, F.P., J. Lu, and S.L. Salzberg, *A review of methods and databases for metagenomic*
672 *classification and assembly*. Brief Bioinform, 2019. **20**(4): p. 1125-1136.
- 673 2. Liu, B., et al., *Rapid Succession of Actively Transcribing Denitrifier Populations in Agricultural Soil*
674 *During an Anoxic Spell*. Front Microbiol, 2018. **9**: p. 3208.
- 675 3. Hanssen, E.N., et al., *Optimizing body fluid recognition from microbial taxonomic profiles*.
676 *Forensic Sci Int Genet*, 2018. **37**: p. 13-20.
- 677 4. Triado-Margarit, X., et al., *Bioaerosols in the Barcelona subway system*. Indoor Air, 2017. **27**(3):
678 p. 564-575.
- 679 5. Segata, N., *On the Road to Strain-Resolved Comparative Metagenomics*. mSystems, 2018. **3**(2).
- 680 6. Zeevi, D., et al., *Structural variation in the gut microbiome associates with host health*. Nature,
681 2019. **568**(7750): p. 43-48.
- 682 7. Liu, M.Y., et al., *Evaluation of ddRADseq for reduced representation metagenome sequencing*.
683 *PeerJ*, 2017. **5**: p. e3837.
- 684 8. Ravi, A., et al., *Comparison of reduced metagenome and 16S rRNA gene sequencing for*
685 *determination of genetic diversity and mother-child overlap of the gut associated microbiota*. J
686 *Microbiol Methods*, 2018. **149**: p. 44-52.
- 687 9. Hess, M.K., et al., *A restriction enzyme reduced representation sequencing approach for low-*
688 *cost, high-throughput metagenome profiling*. PLoS One, 2020. **15**(4): p. e0219882.
- 689 10. Vos, P., et al., *AFLP: a new technique for DNA fingerprinting*. Nucleic Acids Res, 1995. **23**(21): p.
690 4407-14.
- 691 11. Lowry, D.B., et al., *Breaking RAD: an evaluation of the utility of restriction site-associated DNA*
692 *sequencing for genome scans of adaptation*. Mol Ecol Resour, 2017. **17**(2): p. 142-152.
- 693 12. Vendrami, D.L.J., J. Forcada, and J.I. Hoffman, *Experimental validation of in silico predicted RAD*
694 *locus frequencies using genomic resources and short read data from a model marine mammal*.
695 *BMC Genomics*, 2019. **20**(1): p. 72.
- 696 13. Edgar, R.C., *MUSCLE: multiple sequence alignment with high accuracy and high throughput*.
697 *Nucleic Acids Res*, 2004. **32**(5): p. 1792-7.
- 698 14. Popescu, A.A., K.T. Huber, and E. Paradis, *ape 3.0: New tools for distance-based phylogenetics*
699 *and evolutionary analysis in R*. Bioinformatics, 2012. **28**(11): p. 1536-7.
- 700 15. Jain, C., et al., *High throughput ANI analysis of 90K prokaryotic genomes reveals clear species*
701 *boundaries*. Nat Commun, 2018. **9**(1): p. 5114.
- 702 16. DaCosta, J.M. and M.D. Sorenson, *Amplification biases and consistent recovery of loci in a*
703 *double-digest RAD-seq protocol*. PLoS One, 2014. **9**(9): p. e106713.

- 704 17. Wood, D.E., J. Lu, and B. Langmead, *Improved metagenomic analysis with Kraken 2*. Genome Biol, 2019. **20**(1): p. 257.
- 705
- 706 18. Meyer, F., et al., *Assessing taxonomic metagenome profilers with OPAL*. Genome Biol, 2019.
- 707 **20**(1): p. 51.
- 708 19. Ondov, B.D., et al., *Mash: fast genome and metagenome distance estimation using MinHash*.
- 709 Genome Biol, 2016. **17**(1): p. 132.
- 710 20. Segata, N., et al., *Metagenomic microbial community profiling using unique clade-specific*
- 711 *marker genes*. Nature Methods, 2012. **9**(8): p. 811-814.
- 712 21. Mitchell, A.L., et al., *MGNify: the microbiome analysis resource in 2020*. Nucleic Acids Res, 2020.
- 713 **48**(D1): p. D570-D578.
- 714 22. Gloor, G.B., et al., *It's all relative: analyzing microbiome data as compositions*. Ann Epidemiol,
- 715 2016. **26**(5): p. 322-9.
- 716 23. Hippe, B., et al., *Faecalibacterium prausnitzii phylotypes in type two diabetic, obese, and lean*
- 717 *control subjects*. Benef Microbes, 2016. **7**(4): p. 511-7.
- 718 24. Lopez-Siles, M., et al., *Faecalibacterium prausnitzii: from microbiology to diagnostics and*
- 719 *prognostics*. ISME J, 2017. **11**(4): p. 841-852.
- 720 25. Truong, D.T., et al., *Microbial strain-level population structure and genetic diversity from*
- 721 *metagenomes*. Genome Res, 2017. **27**(4): p. 626-638.
- 722 26. Schaeffer, L., et al., *Pseudoalignment for metagenomic read assignment*. Bioinformatics, 2017.
- 723 **33**(14): p. 2082-2088.
- 724 27. Lu, J., et al., *Bracken: estimating species abundance in metagenomics data*. PeerJ Computer
- 725 Science, 2017. **3**.
- 726 28. Almeida, A., et al., *A new genomic blueprint of the human gut microbiota*. Nature, 2019.
- 727 **568**(7753): p. 499-504.
- 728 29. Forster, S.C., et al., *A human gut bacterial genome and culture collection for improved*
- 729 *metagenomic analyses*. Nat Biotechnol, 2019. **37**(2): p. 186-192.
- 730 30. Nayfach, S., et al., *New insights from uncultivated genomes of the global human gut*
- 731 *microbiome*. Nature, 2019. **568**(7753): p. 505-510.
- 732 31. Pasolli, E., et al., *Extensive Unexplored Human Microbiome Diversity Revealed by Over 150,000*
- 733 *Genomes from Metagenomes Spanning Age, Geography, and Lifestyle*. Cell, 2019. **176**(3): p. 649-
- 734 662 e20.
- 735 32. Zou, Y., et al., *1,520 reference genomes from cultivated human gut bacteria enable functional*
- 736 *microbiome analyses*. Nat Biotechnol, 2019. **37**(2): p. 179-185.
- 737 33. Pearce, M.M., et al., *The female urinary microbiome: a comparison of women with and without*
- 738 *urgency urinary incontinence*. mBio, 2014. **5**(4): p. e01283-14.
- 739 34. Stoddard, S.F., et al., *rrnDB: improved tools for interpreting rRNA gene abundance in bacteria*
- 740 *and archaea and a new foundation for future development*. Nucleic Acids Res, 2015.
- 741 **43**(Database issue): p. D593-8.
- 742 35. Ounit, R., et al., *CLARK: fast and accurate classification of metagenomic and genomic sequences*
- 743 *using discriminative k-mers*. BMC Genomics, 2015. **16**: p. 236.
- 744 36. Truong, D.T., et al., *MetaPhlan2 for enhanced metagenomic taxonomic profiling*. Nat Methods,
- 745 2015. **12**(10): p. 902-3.
- 746 37. Menzel, P., K.L. Ng, and A. Krogh, *Fast and sensitive taxonomic classification for metagenomics*
- 747 *with Kaiju*. Nat Commun, 2016. **7**: p. 11257.
- 748 38. Rognes, T., et al., *VSEARCH: a versatile open source tool for metagenomics*. PeerJ, 2016. **4**: p.
- 749 e2584.
- 750 39. Huang, W., et al., *ART: a next-generation sequencing read simulator*. Bioinformatics, 2012. **28**(4):
- 751 p. 593-4.

- 752 40. Lindgreen, S., K.L. Adair, and P.P. Gardner, *An evaluation of the accuracy and speed of*
753 *metagenome analysis tools*. *Sci Rep*, 2016. **6**: p. 19233.
- 754 41. McIntyre, A.B.R., et al., *Comprehensive benchmarking and ensemble approaches for*
755 *metagenomic classifiers*. *Genome Biology*, 2017. **18**(1).
- 756 42. Sczyrba, A., et al., *Critical Assessment of Metagenome Interpretation-a benchmark of*
757 *metagenomics software*. *Nat Methods*, 2017. **14**(11): p. 1063-1071.

758

759

760

761

762 **Figures**

763 Figure 1. An illustration of the suggested RMS profiling procedure. The left branch is executed once for a
764 collection of reference genomes, except that the clustering of genomes may be done at various
765 resolutions depending on later use. The right branch is done for each set of samples. The supplied R
766 package has a tutorial illustrating all steps.

767

768 Figure 2. The distribution of RMS fragment number and lengths in some selected genomes, using the
769 restriction cut sites GAATTC (EcoRI) and TTAA (MseI). There are 10 genomes for each of the species. In
770 the upper panel each dot corresponds to a genome, and the grey sector indicates where the dots should
771 be if it was purely random DNA. In the lower panel each density is based on data from 10 genomes.

772

773 Figure 3. The dendrograms display hierarchical clustering of the 27 genomes in the mock study. In panel
774 A the distances are p-distances computed from a multiple alignment of the 16S sequences from each
775 genome. In panel B the p-distances are based on whole-genome comparisons, and in panel C we used
776 correlation distances based on RMS fragment copy numbers.

777

778 Figure 4. In panel A fragment signal (relative read counts) is plotted against fragment length, and in
779 panel B against fragment GC-content. Each dot corresponds to a fragment cluster, and data from the
780 four single-genome samples are displayed together. In panel C is shown the effect of the simple length-
781 normalization on a single sample. Raw read counts are normalized as described in the text. The red-
782 brown color highlights the fragments within the length-interval 30-500 bases. Note the log-transformed
783 y-axes in all panels.

784

785 Figure 5. Actual and estimated composition of the mock with 20 genomes and no extra genomes spiked-
786 in. The extra strains, absent from the samples, are indicated by the tan colors. The left bar is the actual
787 compositions, and the three additional bars (A, B and C) show estimates from the three replicates of this
788 sample.

789

790 Figure 6. Each panel show actual versus estimated proportions for all 27 genomes as markers. Each
791 marker color corresponds to a genome, and the three marker types are the three replicates. Species
792 with a single genome is colored in shades of grey, while those where two or three genomes are similar
793 have more distinct coloring. There are eight different mocks, all consisting of the 20 genomes in the
794 original mock, but with various additional genomes spiked-in. The header above each panel indicate
795 which genome has been spiked-in. The gray line in the background is where marker should be if the
796 estimates were perfect. Note that in the lower left corner (those who are absent and predicted to be
797 absent) many markers of different colors overlap each other.

798

799 Figure 7. The scatter plots show actual versus predicted relative abundances for the simulated data.
800 Each dot is a relative abundance of a genome, and each panel contain results from 25 samples. The
801 upper panels are RMS data estimated by COLS, and the lower panels shotgun data estimated by
802 Kraken2. The resolution of the communities increases from left to right, as indicated by the upper panel
803 headers. The average Manhattan distance D is displayed within each panel.

804

805 Figure 8. In panel A each dot indicates distance to the nearest neighbor from an *E. coli* genome,
806 measured either as RMS correlation distance (x-axis) or whole-genomes p-distance (y-axis). The grey
807 dots are the results for all 1066 genomes, and the black dots are for the 54 genomes left after clustering
808 with maximum condition value 100. In panel B the neighbor joining tree is based on the p-distances
809 between all strains, and the black tips indicate the cluster centroid genomes.

810

811 Figure 9. The left panel A shows a tree for the 54 *E. coli* strains, based on the whole-genome p-distances,
812 with a scale marker in the lower right corner. The leaf nodes are marked by how they were classified as
813 present in a single sample. No leaf marker means classified as absent by both shotgun+Kraken2 and
814 RMS+COLS. Grey markers indicate classified as present by shotgun+Kraken2 only. Black dots indicate
815 classified as present by both methods. No genomes were classified as present by RMS+COLS only. In the
816 right panel B the boxplot shows the number of genomes, out of the 54, estimated to be present by
817 RMS+COLS over time in all infant gut samples (Meconium is newborn feces). There are 94 samples
818 behind each box.

819

820 Figure 10. Re-analysis of the data from [7]. In the left panel A we compare the profiles obtained by using
821 Kraken2 on the conventional shotgun (shotgun+kraken2) data and the RMS data (ddRADseq+kraken2).
822 In the original paper the same comparison was made using MetaPhlan instead of Kraken2. In addition,
823 we also used our approach described in this paper on the RMS data (ddRADseq+COLS). The marker-type
824 indicates the samples, and the coloring the methods. In the right panel B we focus only on a strain
825 resolution of the species *F. prausnitzii*, being the most dominant species in these samples. The nine
826 strains listed are from the MGnify database (<https://www.ebi.ac.uk/metagenomics/>).

827

828

829

830 **Tables**

831 Table 1. For each genome is listed its size (megabasepairs), GC-content and the number of RMS
 832 fragments in the 30-500 length interval. Genomes with an asterisk (*) after its name were spiked-in, and
 833 not part of the original mock.

Genome	Size	GC	RMS-fragments
<i>Acinetobacter baumannii</i> strain 5377	3.98	0.39	961
<i>Schaalia odontolytica</i> strain 1A.21	2.39	0.65	92
<i>Bacillus cereus</i> strain NRS 248	5.22	0.36	2025
<i>Bacteroides vulgatus</i> strain ATCC 8482	5.16	0.42	2047
<i>Clostridium beijerinckii</i> strain NCIMB 8052	6.00	0.30	2463
<i>Cutibacterium acnes</i> strain KPA171202	2.56	0.60	300
<i>Deinococcus radiodurans</i> strain R1	3.06	0.67	115
<i>Enterococcus faecalis</i> ATCC 19433*	2.87	0.38	902
<i>Enterococcus faecalis</i> ATCC 29212*	3.01	0.37	922
<i>Enterococcus faecalis</i> strain OG1RF	2.74	0.38	822
<i>Escherichia coli</i> strain K12 substrain MG1655	4.64	0.51	920
<i>Helicobacter pylori</i> NCTC 11637*	1.60	0.39	233
<i>Helicobacter pylori</i> strain 26695	1.67	0.39	255
<i>Lactobacillus gasseri</i> ATCC 33323*	1.82	0.35	662
<i>Lactobacillus gasseri</i> strain 63 AM	1.89	0.35	675
<i>Listeria monocytogenes</i> strain EGDe	2.94	0.38	1504
<i>Neisseria meningitidis</i> strain MC58	2.27	0.52	332
<i>Pseudomonas aeruginosa</i> strain PAO1-LAC	6.26	0.66	143
<i>Rhodobacter sphaeroides</i> strain ATH 2 4 1	4.13	0.69	92
<i>Staphylococcus aureus</i> strain TCH1516	2.88	0.33	861
<i>Staphylococcus epidermidis</i> FDA strain PCI 1200	2.50	0.32	854
<i>Streptococcus agalactiae</i> ATCC 13813*	2.11	0.35	661
<i>Streptococcus agalactiae</i> strain 2603 V R	2.16	0.36	692
<i>Streptococcus mutans</i> ATCC 25175*	1.99	0.37	671
<i>Streptococcus mutans</i> strain UA159	2.03	0.37	680
<i>Streptococcus pneumoniae</i> ATCC 6305*	2.02	0.40	709
<i>Streptococcus pneumoniae</i> strain TIGR4	2.16	0.40	771

834

835

836 **Additional Files**837 **Additional file 1 --- Supplementary figure 1**

838 All simulated samples contained reads from 100 randomly selected genomes, and their relative
 839 abundances in the sample were according to this barplot. The largest abundance is 100 times the
 840 smallest. Different genomes were selected as the most/least abundant and absent ones in each sample,
 841 but this abundance distribution was used every time.

842

843 Additional file 2 --- Supplementary figure 2

844 In order to simulate RMS data, some known biases were introduced to the signals. The upper panel
845 shows the fragment-length bias used. All signals were scaled by this function, i.e. fragments of length
846 around 200 bases remained close to unchanged (scale \$1.0\$) while signals from shorter or longer
847 fragments were scaled down. The lower panel shows the fragment-bias distribution. For each fragment
848 within a genome, a factor was sampled from this distribution, and the signals from the fragments were
849 scaled accordingly. The mean value of this distribution is \$1.0\$, but some fragments may have signals up
850 to six times as large, or down to almost nothing. Both the length-bias function and the fragment-bias
851 distribution were estimated from real RMS data.

852

853

Figures

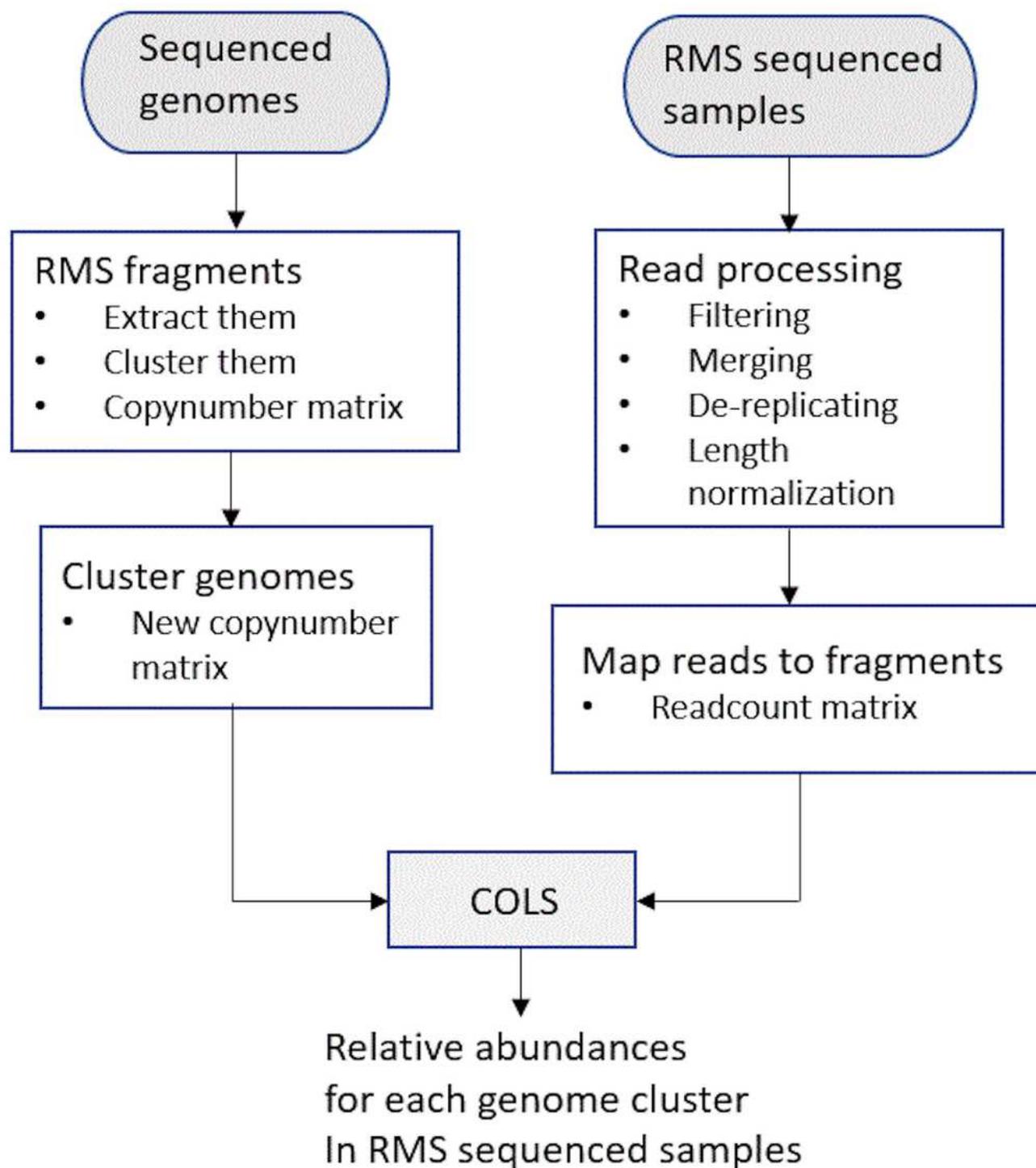


Figure 1

An illustration of the suggested RMS profiling procedure. The left branch is executed once for a collection of reference genomes, except that the clustering of genomes may be done at various resolutions

depending on later use. The right branch is done for each set of samples. The supplied R package has a tutorial illustrating all steps.

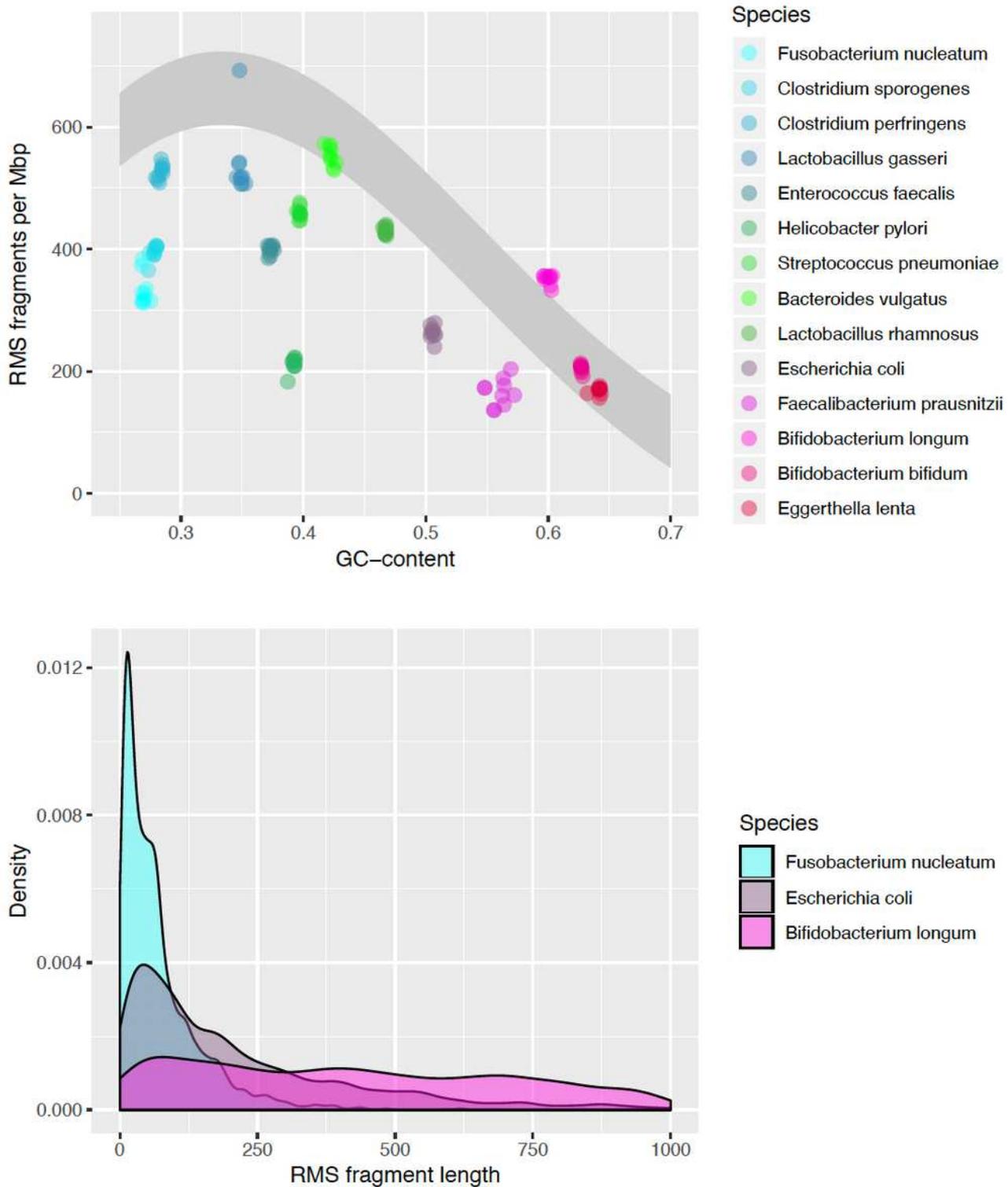


Figure 2

The distribution of RMS fragment number and lengths in some selected genomes, using the restriction cut sites GAATTC (EcoRI) and TTAA (MseI). There are 10 genomes for each of the species. In the upper

panel each dot corresponds to a genome, and the grey sector indicates where the dots should be if it was purely random DNA. In the lower panel each density is based on data from 10 genomes.

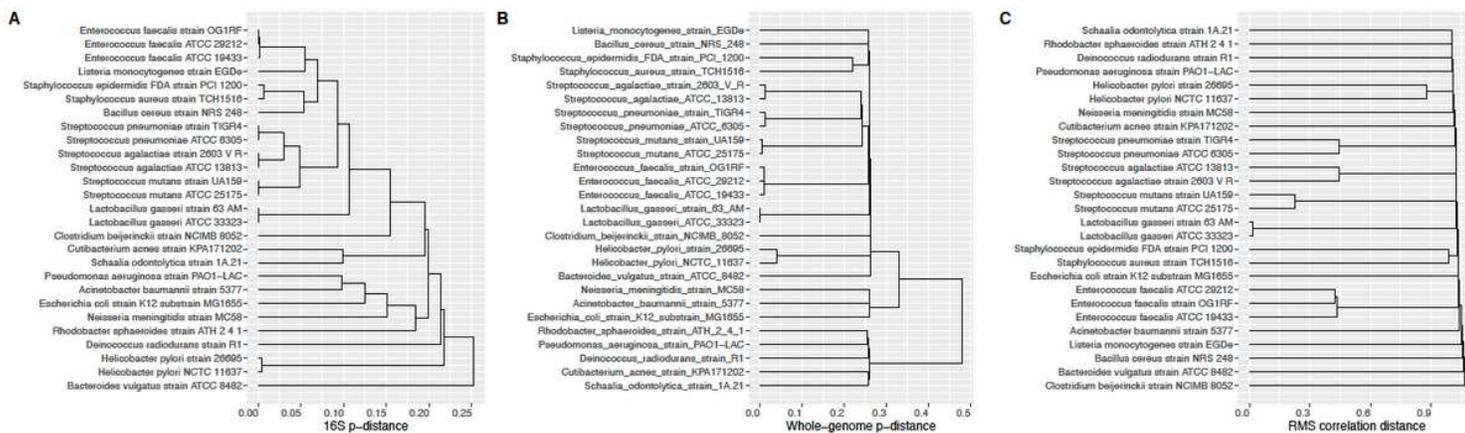


Figure 3

The dendrograms display hierarchical clustering of the 27 genomes in the mock study. In panel A the distances are p-distances computed from a multiple alignment of the 16S sequences from each genome. In panel B the p-distances are based on whole-genome comparisons, and in panel C we used correlation distances based on RMS fragment copy numbers.

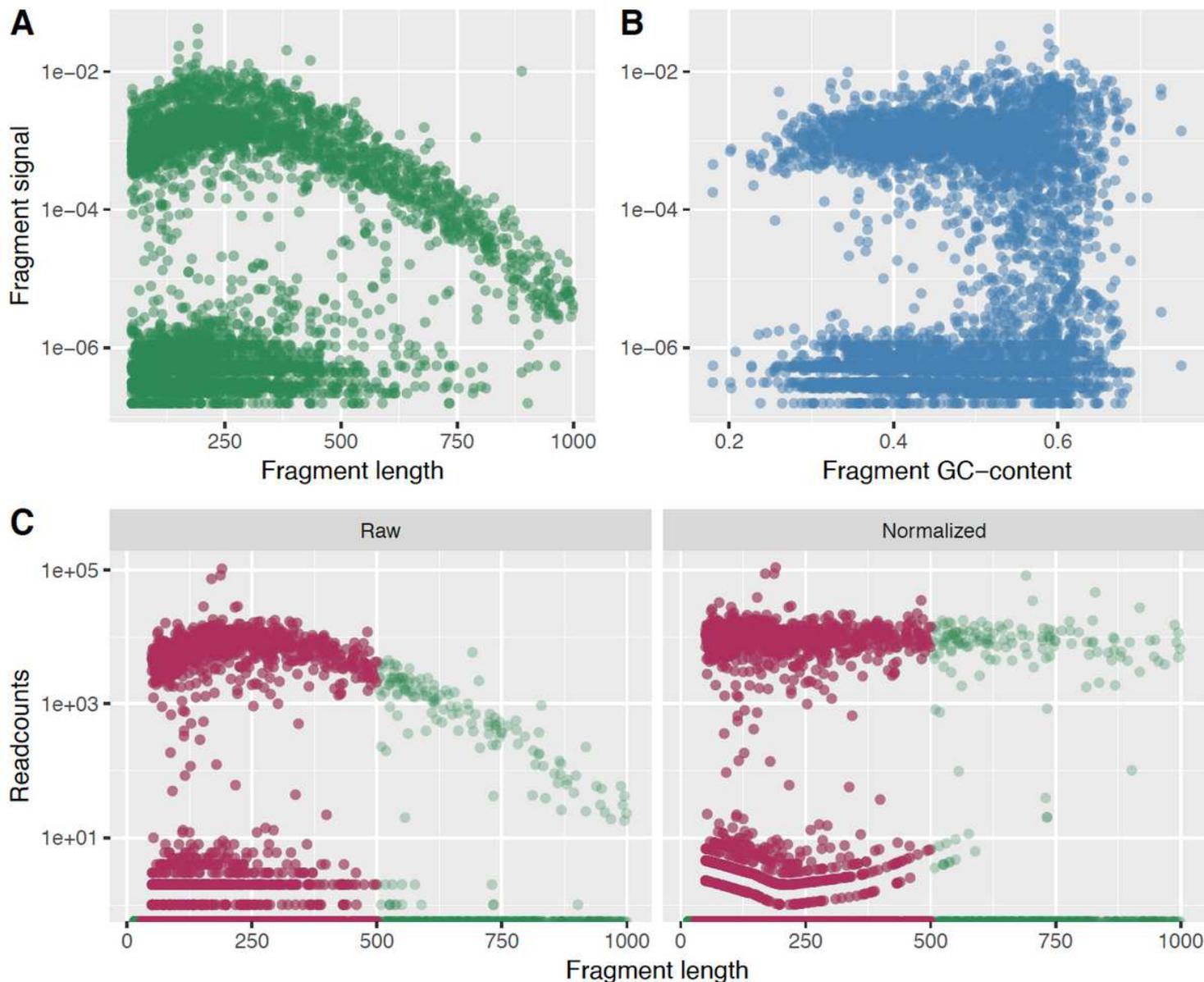


Figure 4

In panel A fragment signal (relative read counts) is plotted against fragment length, and in panel B against fragment GC-content. Each dot corresponds to a fragment cluster, and data from the four single-genome samples are displayed together. In panel C is shown the effect of the simple length-normalization on a single sample. Raw read counts are normalized as described in the text. The red-brown color highlights the fragments within the length-interval 30-500 bases. Note the log-transformed y-axes in all panels.

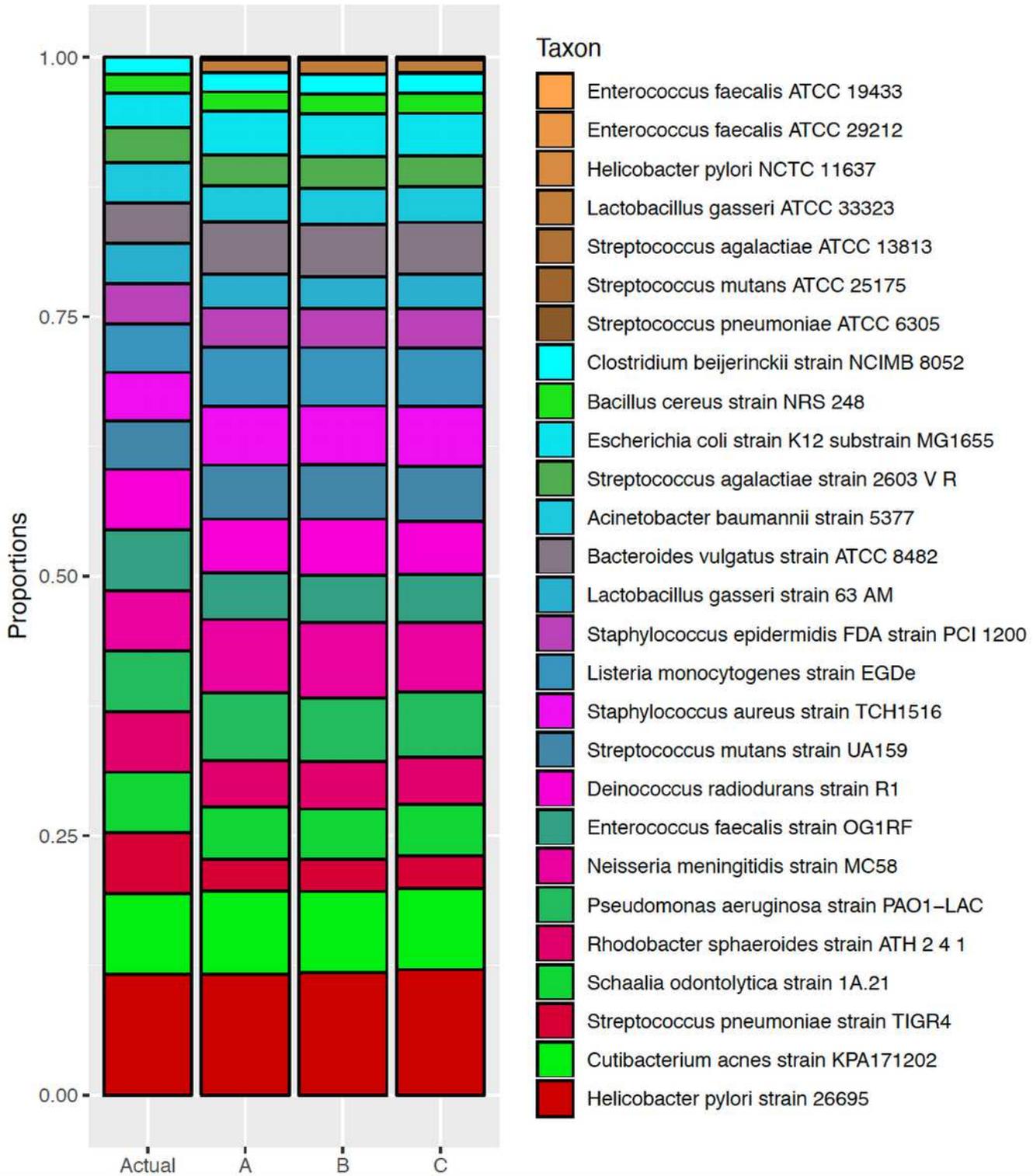


Figure 5

Actual and estimated composition of the mock with 20 genomes and no extra genomes spiked-in. The extra strains, absent from the samples, are indicated by the tan colors. The left bar is the actual compositions, and the three additional bars (A, B and C) show estimates from the three replicates of this sample.

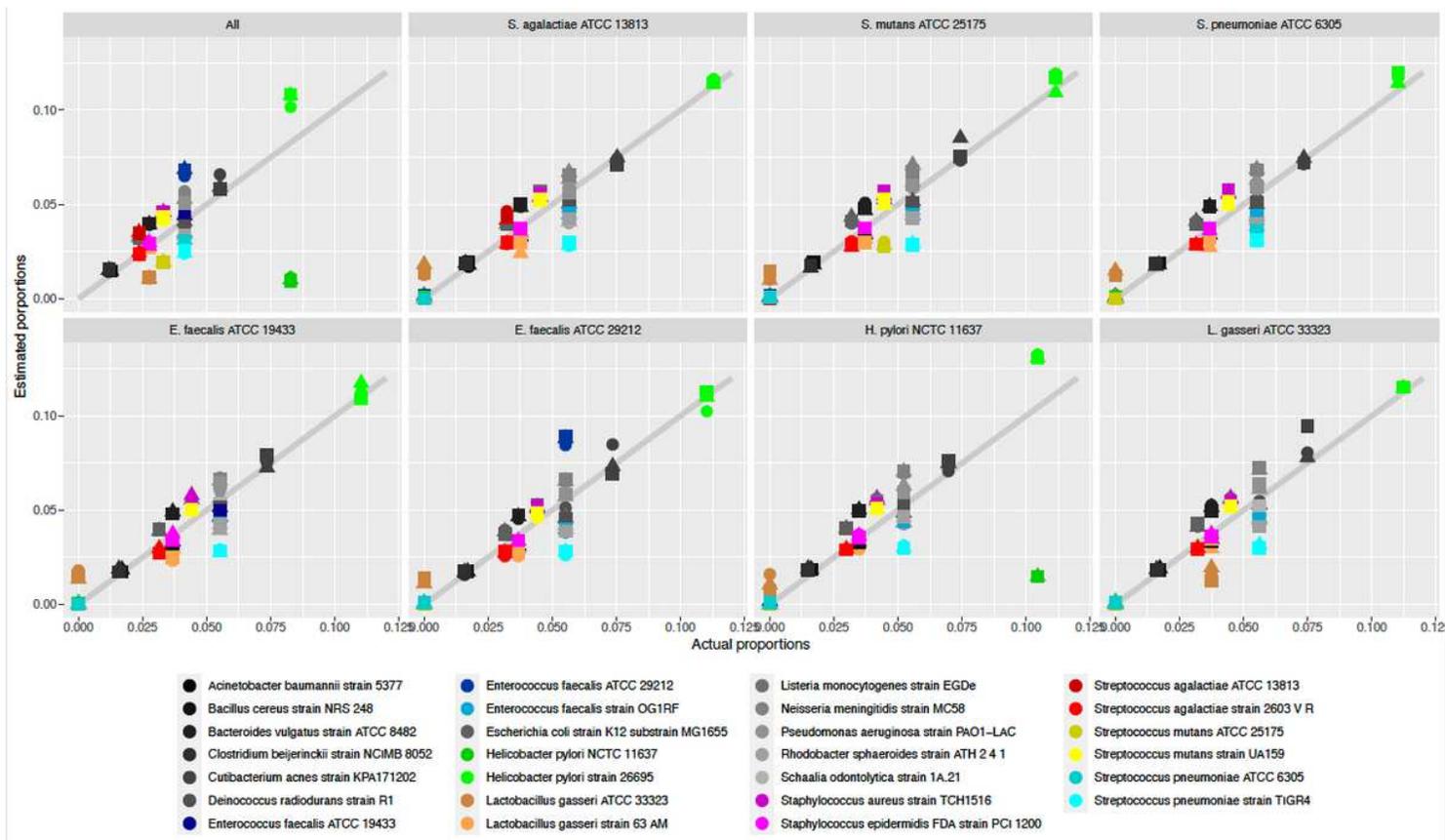


Figure 6

Each panel show actual versus estimated proportions for all 27 genomes as markers. Each marker color corresponds to a genome, and the three marker types are the three replicates. Species with a single genome is colored in shades of grey, while those where two or three genomes are similar have more distinct coloring. There are eight different mocks, all consisting of the 20 genomes in the original mock, but with various additional genomes spiked-in. The header above each panel indicate which genome has been spiked-in. The gray line in the background is where marker should be if the estimates were perfect. Note that in the lower left corner (those who are absent and predicted to be absent) many markers of different colors overlap each other.

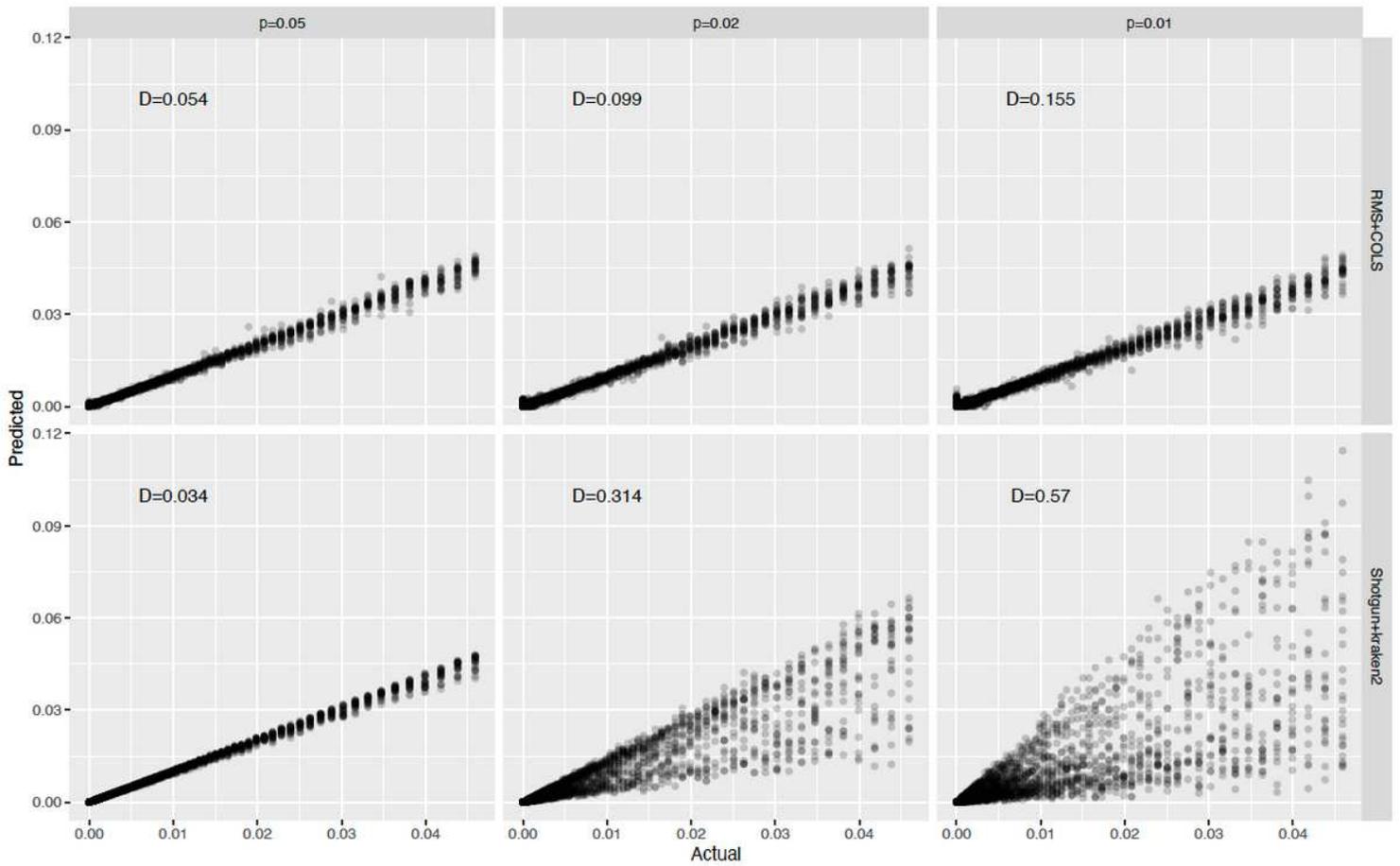


Figure 7

The scatter plots show actual versus predicted relative abundances for the simulated data. Each dot is a relative abundance of a genome, and each panel contain results from 25 samples. The upper panels are RMS data estimated by COLS, and the lower panels shotgun data estimated by Kraken2. The resolution of the communities increases from left to right, as indicated by the upper panel headers. The average Manhattan distance D is displayed within each panel.

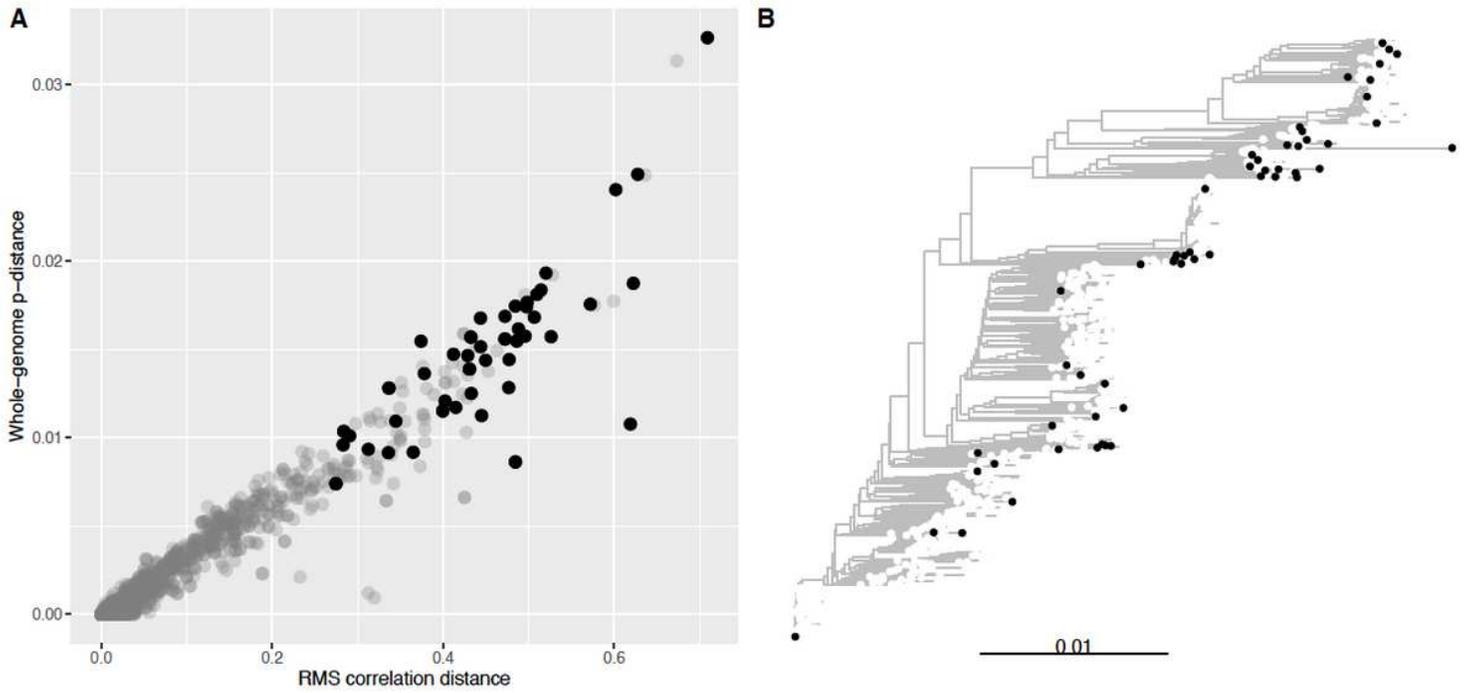


Figure 8

In panel A each dot indicates distance to the nearest neighbor from an *E. coli* genome, measured either as RMS correlation distance (x-axis) or whole-genomes p-distance (y-axis). The grey dots are the results for all 1066 genomes, and the black dots are for the 54 genomes left after clustering with maximum condition value 100. In panel B the neighbor joining tree is based on the p-distances between all strains, and the black tips indicate the cluster centroid genomes.

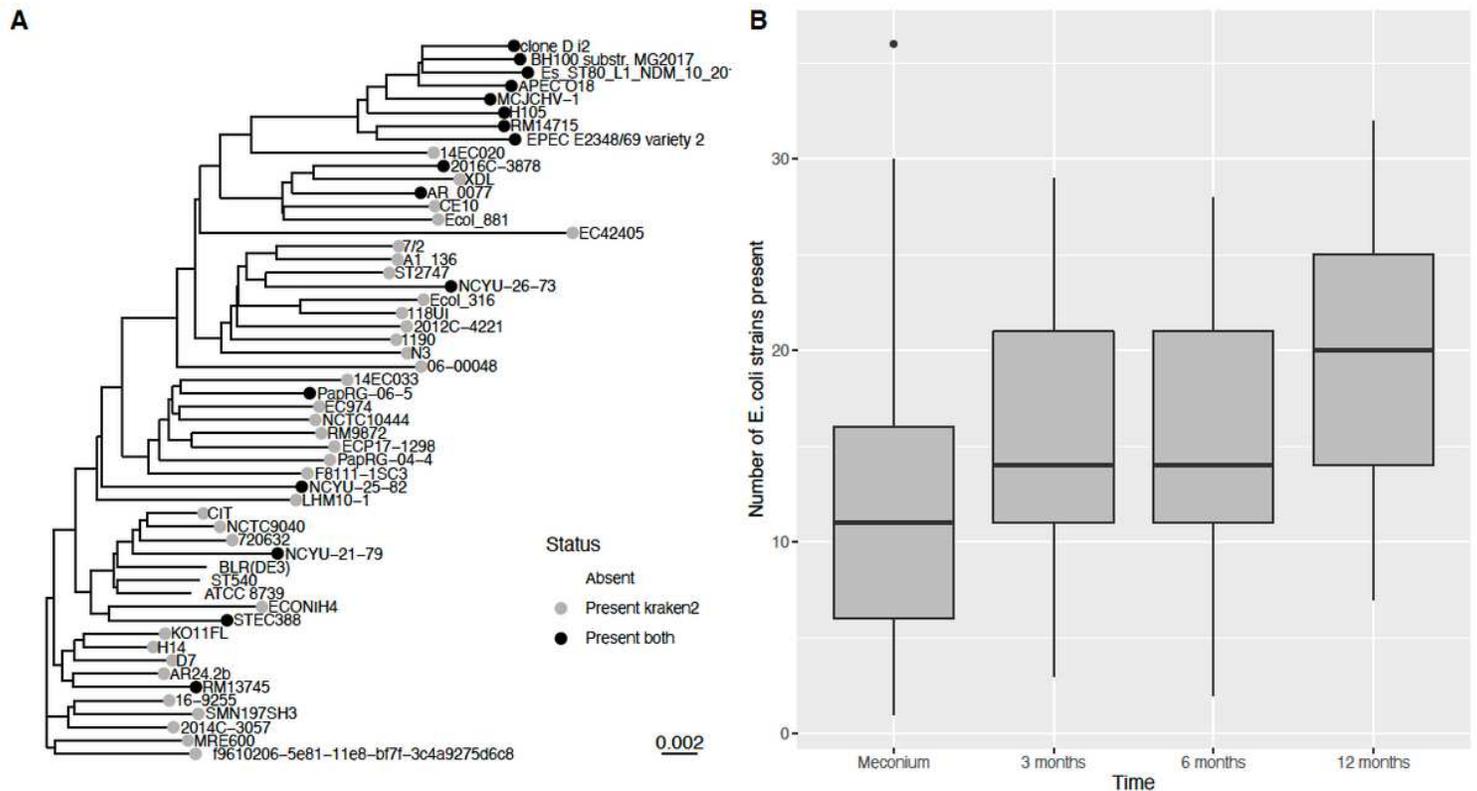


Figure 9

The left panel A shows a tree for the 54 *E. coli* strains, based on the whole-genome p-distances, with a scale marker in the lower right corner. The leaf nodes are marked by how they were classified as present in a single sample. No leaf marker means classified as absent by both shotgun+Kraken2 and RMS+COLS. Grey markers indicate classified as present by shotgun+Kraken2 only. Black dots indicate classified as present by both methods. No genomes were classified as present by RMS+COLS only. In the right panel B the boxplot shows the number of genomes, out of the 54, estimated to be present by RMS+COLS over time in all infant gut samples (Meconium is newborn feces). There are 94 samples behind each box.

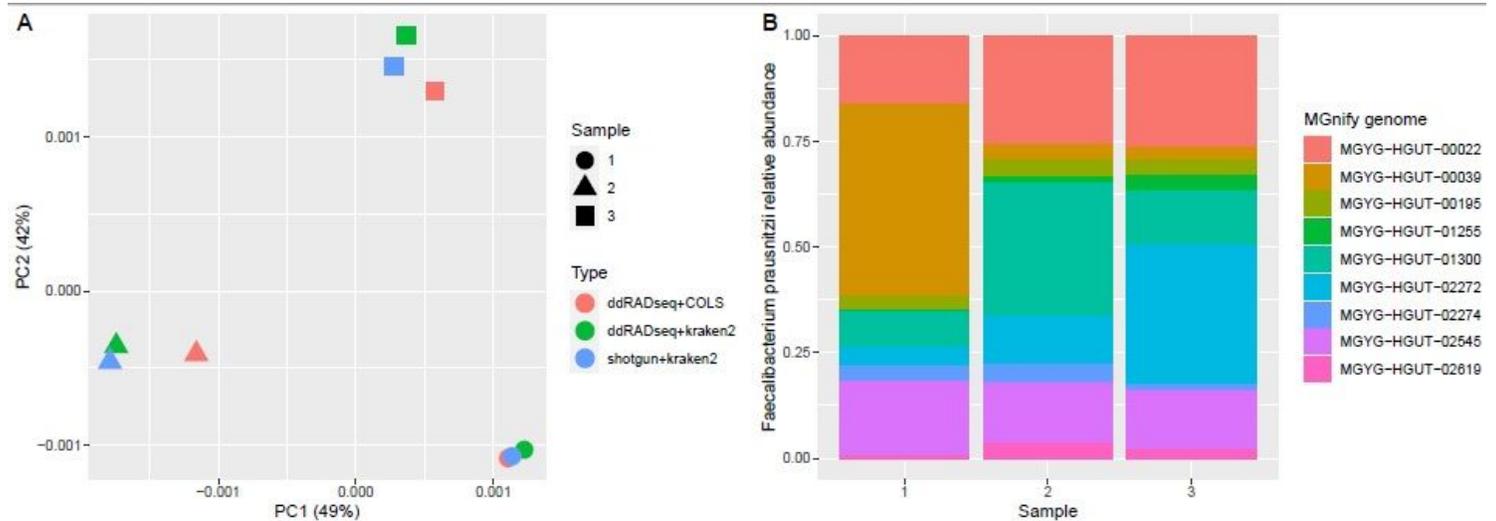


Figure 10

Re-analysis of the data from [7]. In the left panel A we compare the profiles obtained by using Kraken2 on the conventional shotgun (shotgun+kraken2) data and the RMS data (ddRADseq+kraken2). In the original paper the same comparison was made using MetaPhlan instead of Kraken2. In addition, we also used our approach described in this paper on the RMS data (ddRADseq+COLS). The marker-type indicates the samples, and the coloring the methods. In the right panel B we focus only on a strain resolution of the species *F. prausnitzii*, being the most dominant species in these samples. The nine strains listed are from the MGnify database (<https://www.ebi.ac.uk/metagenomics/>).

Supplementary Files

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

- [supfig1.pdf](#)
- [supfig2.pdf](#)