

1 **Monitoring the microbiome for food safety and quality using deep** 2 **shotgun sequencing**

3

4 **Running Title:** Monitoring the microbiome for food safety and quality

5

6 Kristen L. Beck^{1,8*+}, Niina Haiminen^{2,8+}, David Chambliss^{1,8}, Stefan Edlund^{1,8}, Mark Kunitomi^{1,8},
7 B. Carol Huang^{3,8}, Nguyet Kong^{3,8}, Balasubramanian Ganesan^{4,5,8}, Robert Baker^{4,8}, Peter
8 Markwell^{4,8}, Ban Kawas^{1,8}, Matthew Davis^{1,8}, Robert J. Prill^{1,8}, Harsha Krishnareddy^{1,8}, Ed
9 Seabolt^{1,8}, Carl H. Marlowe^{6,8}, Sophie Pierre^{7,8}, André Quintanar^{7,8}, Laxmi Parida^{2,8}, Geraud
10 Dubois^{1,8}, James Kaufman^{1,8}, and Bart C. Weimer^{3,8*}

11

12 *Co-corresponding authors

13 +Contributed equally

14

15 Contact information: Kristen L. Beck, IBM Almaden Research Center, 650 Harry Road, San Jose
16 CA, 95120 USA, klbeck@us.ibm.com, +1 408-927-1963

17

18 **AUTHOR AFFILIATIONS:**

19 ¹IBM Almaden Research Center, San Jose CA

20 ²IBM T.J. Watson Research Center, Yorktown Heights, NY

21 ³University of California Davis, School of Veterinary Medicine, 100K Pathogen Genome Project,
22 Davis, CA 95616

23 ⁴Mars Global Food Safety Center, Beijing, China

24 ⁵Wisdom Health, A Division of Mars Petcare, Vancouver WA.

25 ⁶Bio-Rad Laboratories, Hercules CA

26 ⁷Bio-Rad, Food Science Division, MArnes-La-Coquette, France

27 ⁸Consortium for Sequencing the Food Supply Chain, San Jose, CA

28

29

30 **ABSTRACT:**

31 In this work, we hypothesized that shifts in the food microbiome can be used as an indicator of
32 unexpected contaminants or environmental changes. To test this hypothesis, we sequenced total
33 RNA of 31 high protein powder (HPP) samples of poultry meal pet food ingredients. We
34 developed a microbiome analysis pipeline employing a key eukaryotic matrix filtering step that
35 improved microbe detection specificity to >99.96% during *in silico* validation. The pipeline
36 identified 119 microbial genera per HPP sample on average with 65 genera present in all
37 samples. The most abundant of these were *Bacteroides*, *Clostridium*, *Lactococcus*, *Aeromonas*,
38 and *Citrobacter*. We also observed shifts in the microbial community corresponding to
39 ingredient composition differences. When comparing culture-based results for *Salmonella* with
40 total RNA sequencing, we found that *Salmonella* growth did not correlate with multiple
41 sequence analyses. We conclude that microbiome sequencing is useful to characterize complex
42 food microbial communities, while additional work is required for predicting specific species'
43 viability from total RNA sequencing.

44

45 **KEYWORDS:**

46 microbiome, food safety, bioinformatics, shotgun sequencing, microbial ecology, pathogens

47

48

49 **1. INTRODUCTION:**

50 Sequencing the microbiome of food may reveal characteristics about the associated
51 microbial content that culturing or targeted whole genome sequencing alone cannot. However, to
52 meet the various needs of food safety and quality, next generation sequencing (NGS) and analysis
53 techniques require additional development¹ with specific consideration for accuracy, speed, and
54 applicability across the supply chain.² Microbial communities and their characteristics have been
55 studied in relation to flavor and quality in fermented foods,³⁻⁵ agricultural processes in grape⁶ and
56 apple fruit⁷, and manufacturing processes and production batches in Cheddar cheese.⁸ However,
57 the advantage of using the microbiome specifically for food safety and quality has yet to be
58 demonstrated.

59 Currently, food safety regulatory agencies including the Food and Drug Administration
60 (FDA), Centers for Disease Control and Prevention (CDC), United States Department of
61 Agriculture (USDA), and European Food Safety Authority (EFSA) are converging on the use of
62 whole genome sequencing (WGS) for pathogen detection and outbreak investigation. Large scale
63 WGS of food-associated bacteria was first initiated via the 100K Pathogen Genome Project⁹ with
64 the goal of expanding the diversity of bacterial reference genomes— a crucial need for foodborne
65 illness outbreak investigation, traceability, and microbiome studies.^{10,11} However, since WGS
66 relies on culturing a microbial isolate prior to sequencing, there are inherent biases and limitations
67 in its ability to describe the microorganisms and their interactions in a food sample. Such
68 information would be very valuable for food safety and quality applications.

69 High throughput sequencing of total DNA and total RNA are promising approaches to
70 characterize microbial niches in their native state without introducing bias due to culturing.¹²⁻¹⁴
71 Additionally, total RNA sequencing has the potential to provide evidence of live and biologically
72 active components of the sample.^{14,15} It also provides accurate microbial naming, relative

73 microbial abundance, and better reproducibility than total DNA or amplicon sequencing.¹⁴ Total
74 RNA sequencing minimizes PCR amplification bias that occurs in single gene amplicon
75 sequencing and overcomes the decreased detection sensitivity from using DNA sequencing in
76 metagenomics.¹⁴ Total RNA metatranscriptome sequencing, however, is yet to be examined in raw
77 food ingredients as a method to provide robust characterization of the microbial communities and
78 the interacting population dynamics.

79 From a single sequenced food microbiome, numerous dimensions of the sample can be
80 characterized that may yield important indicators of safety and quality. Using total DNA or RNA,
81 evidence for the eukaryotic food matrix can be examined. In Haiminen *et al.*,¹⁶ we quantitatively
82 demonstrated the utility of metagenome sequencing to authenticate the composition of complex
83 food matrices. In addition, from total DNA or RNA, one can observe signatures from commensal
84 microbes, pathogenic microbes, and genetic information for functional potential (from DNA) or
85 biologically active function (from RNA).^{14,15} Detecting active transcription from live microbes in
86 food is very important to avoid spurious microbial observations that may instead be false positives
87 due to quiescent DNA in the sample. Use of RNA in food analytics also offers the opportunity to
88 examine expression of metabolic processes that are related to antibiotic resistance,^{17,18} virulence
89 factors, or replication genes, among others. Additionally, it has the potential to define viable
90 microbes that are capable of replication in the food and even microorganisms that stop replicating
91 but continue to produce metabolic activity that changes food quality and safety.¹⁹⁻²⁴

92 Microorganisms are sensitive to changes in temperature, salinity, pH, oxygen content, and
93 many other physicochemical factors that alter their ability to grow, persist, and cause disease. They
94 exist in dynamic communities that change in response to environmental perturbation – just as the
95 gut microbiome shifts in response to diet.²⁵⁻²⁸ Shifts in microbiome composition or activity can be
96 leveraged in the application of microbiome characterization to monitor the food supply chain. For

97 example, Noyes et al. followed the microbiome of cattle from the feed lot to the food packaging,
98 concluding that the microbial community and antibiotic resistance characteristics change based on
99 the processing stage.^{17,18,29} We hypothesize that observable shifts in microbial communities of
100 food can serve as an indicator of food quality and safety.

101 In this work, we examined 31 high protein powder samples (HPP; derived from poultry meal).
102 HPP are commonly used raw materials in pet foods. They are subject to microbial growth prior to
103 preparation and continued survival in powder form.³⁰ We subjected the HPP samples to deep total
104 RNA sequencing with ~300 million reads per sample. In order to process the 31 samples collected
105 over ~1.5 years from two suppliers at a single location, we defined and calibrated the appropriate
106 methods– from sample preparation to bioinformatic analysis– needed to taxonomically identify
107 the community members present and to detect key features of microbial growth. First, we removed
108 the HPP’s food matrix RNA content as eukaryotic background with an important bioinformatic
109 filtering step designed specifically for food analysis. The remaining sequences were used for
110 relative quantification of microbiome members and for identifying shifts based on food matrix
111 content, production source, and *Salmonella* culturability. This work demonstrates that total RNA
112 sequencing is a robust approach for monitoring the food microbiome for use in food safety and
113 quality applications, while additional work is required for predicting pathogen viability.

114

115 **2. RESULTS:**

116 **2.1 Evaluation of microbial identification capability in total RNA and DNA sequencing**

117 Microbial identification in microbiomes often leverages shotgun DNA sequencing; however,
118 total RNA sequencing can provide additional information about viable bacterial activity in a
119 community via transcriptional activity. Since using total RNA to study food microbiomes is novel,
120 each step of the analysis workflow (Figure 1) was carefully designed and scrutinized for accuracy.

121 For all analyses done in this study, we report relative abundance in reads per million (RPM)
122 (Equation 1) as recommended by Gloor et al^{31,32} and apply the conservative threshold of RPM >
123 0.1 to indicate presence as previously described by Langelier et al and Illot et al.^{33,34} Numerically,
124 this threshold translates to ~30 reads per genus per sample considering a sequencing depth of ~300
125 million reads per sample (Methods Section 4.4). First, we examined the effectiveness of RNA for
126 taxonomic identification and relative quantification of microbes in the presence of food matrix
127 reads. We observed that RNA sequencing results correlated ($R^2 = 0.93$) with the genus relative
128 quantification provided by DNA sequencing (Supplementary Figure S1). RNA sequencing also
129 detected more genera demonstrated by a higher α -diversity than the use of DNA (Supplementary
130 Figure S2). Additionally, from the same starting material, total RNA sequencing resulted in 2.4-
131 fold more reads classified to microbial genera compared to total DNA sequencing (after
132 normalizing for sequencing depth). This increase is substantial as microbial reads are such a small
133 fraction of the total sequenced reads. Considering these results, we further examined the microbial
134 content from total RNA extracted from 31 high protein powder (HPP) samples (Supplementary
135 Table 1) that resulted in an average of ~300 million paired end 150 bp sequencing reads per sample
136 in this study.

137

138 **2.2 Evaluation and application of *in silico* filtering of eukaryotic food matrix reads**

139 Sequenced reads from the eukaryotic host or food matrix may lead to false positives for microbial
140 identification in microbiome studies.³⁵ This may occur partly due to reads originating from low
141 complexity regions of eukaryotic genomes, e.g. telomeric and centromeric repeats, being
142 misclassified as spurious microbial hits.³⁶ In total DNA or RNA sequencing of clinical or animal
143 or even plant microbiomes, eukaryotic content may often comprise > 90% of the total sequencing
144 reads. This presents an important bioinformatic challenge that we addressed by filtering matrix

145 content using a custom-built reference database of 31 common food ingredient and contaminant
146 genomes (Supplementary Table 2) using the *k*-mer classification tool Kraken.³⁷ This step allows
147 for rapidly classifying all sequenced reads (~300 million reads for each of 31 samples) as matrix
148 or non-matrix. The matrix filtering process yielded an estimate of the total percent matrix content
149 for a sample. See our work in Haiminen et al.³⁸ on quantifying the eukaryotic food matrix
150 components with further precision.

151 To validate the matrix filtering step, we constructed *in silico* mock food microbiomes with
152 a high proportion of complex food matrix content and low microbial content (Supplementary Table
153 3). We then computed the true positive, false positive, and false negative rates of observed
154 microbial genera and sequenced reads (Table 1). False positive viral, archaeal, and eukaryotic
155 microbial genera (as well as bacteria) were observed without matrix filtering, although bacteria
156 were the only microbes included in the simulated mixtures. Introducing a matrix filtering step to
157 the pipeline improved read classification specificity to >99.96% (from 78–93% without filtering)
158 in both simulated food mixtures, while maintaining zero false negatives. With this level of
159 demonstrated accuracy, we used bioinformatic matrix filtering prior to further microbiome
160 analysis.

161

162 **2.3 High protein powder microbiome ecology**

163 After filtering eukaryotic matrix sequences, we applied the remaining steps in the
164 bioinformatic workflow (Figure 1) to examine the shift in the high protein powder (HPP)
165 microbiome membership and to quantify the relative abundance of microbes at the genus level.
166 Genus is the first informative taxonomic rank for food pathogen identification that can be
167 considered accurate given current incompleteness of reference databases^{11,39–42} and was therefore
168 used in subsequent analyses. Overall, between 98 and 195 microbial genera (avg. 119) were

169 identified (RPM > 0.1) per HPP sample (Supplementary Table 4). When analyzing α -diversity
170 i.e. the number of microbes detected per sample, inter-sample comparisons may become skewed
171 unless a common number of reads is considered since deeper sequenced samples may contain more
172 observed genera merely due to a greater sampling depth.^{43,44} Thus, we utilized bioinformatic
173 rarefaction i.e. subsampling analysis to showcase how microbial diversity was altered by
174 sequencing depth. Examination of α -diversity across a range of *in silico* subsampled sequencing
175 depths showed that the community diversity varied across samples (Figure 2A). One sample
176 (MFMB-04) had 1.7 times more genera (195) than the average across other samples (avg. 116,
177 range 98–143) and exhibited higher α -diversity than any other sample at each *in silico* sampled
178 sequencing depth (Figure 2A). Rarefaction analysis further demonstrated that when considering
179 fewer than ~67 million sequenced reads, the observable microbial population was not saturated
180 (median elbow calculated as indicated in Satopää, et al.⁴⁵). This observation suggests that deeper
181 sequencing or more selective sequencing of the HPP microbiomes will reveal more microbial
182 diversity.

183 Notably, between 2%–4% (approximately 5,000,000–14,000,000) of reads per sample
184 remained unclassified as either eukaryotic matrix or microbe (Supplementary Figure S3).
185 However, the unclassified reads exhibited a GC (guanine plus cytosine) distribution similar to
186 reads classified as microbial (Supplementary Figure S4) indicating these reads may represent
187 microbial content that is absent or sufficiently divergent from existing references.

188 We calculated β -diversity to study inter-sample microbiome differences and to identify any
189 potential outliers among the sample collection. The Aitchison distances⁴⁶ of microbial relative
190 abundances were calculated between samples (as recommended for compositional microbiome
191 data^{31,32}), and the samples were hierarchically clustered based on the resulting distances (Figure

192 2B). The two primary clades were mostly defined by supplier (except for MFMB-17). In Haiminen
193 *et al.*,³⁸ we reported that three of the HPP samples contained unexpected eukaryotic species. We
194 hypothesized that the presence of these contaminating matrix components (beef identifiable as *Bos*
195 *taurus* and pork identifiable as *Sus scrofa*) would alter the microbiome as compared to chicken
196 (identifiable as *Gallus gallus*) alone. Clustering HPP samples using their microbiome membership
197 led to a distinctly different group of the matrix-contaminated samples, supporting this hypothesis
198 (Figure 2B). These observations indicate that samples can be discriminated based on their
199 microbiome content for originating source and supplier, which is necessary for source tracking
200 potential hazards in food.

201 **2.4 Comparative analysis of high protein powder microbiome membership and** 202 **composition**

203 We identified 65 genera present in all HPP samples (Figure 3A), whose combined
204 abundance accounted for between 88-99% of the total abundances of detected genera per sample.
205 *Bacteroides*, *Clostridium*, *Lactococcus*, *Aeromonas*, and *Citrobacter* were the five most abundant
206 of these microbial genera. The identified microbial genera also included viruses, the most abundant
207 of which was *Gyrovirus* (< 10 RPM per sample). *Gyrovirus* represents a genus of non-enveloped
208 DNA viruses responsible for chicken anemia which is ubiquitous in poultry. While there were only
209 65 microbial genera identified in all 31 HPP samples, the α -diversity per sample was on average
210 two-fold greater as previously indicated.

211 Beyond the collection of 65 microbes observed in all samples, there were an additional 164
212 microbes present in various HPP samples. Together, we identified a total of 229 genera among the
213 31 HPP samples tested (Figure 3B and 4, Supplementary Table 4). In order to identify genera that
214 were most variable between samples, we computed the median absolute deviation (MAD)⁴⁷ using
215 the normalized relative abundance of each microbe (Figure 5A). The abundance of *Bacteroides*

216 was the most variable among samples (median = 148.1 RPM, MAD = 30.6) and showed increased
217 abundance in almost all samples from Supplier A compared to Supplier B (abundance for the 10
218 most variable genera shown in Figure 5B). *Clostridium* (median = 37.4 RPM, MAD = 24.2),
219 *Lactococcus* (median = 36.8 RPM, MAD = 18.2), and *Lactobacillus* (median = 24.2, MAD = 7.2)
220 were also highly variable and 3–4 fold more abundant in samples MFMB-04 and MFMB-20
221 compared to other samples (Figure 5B). *Pseudomonas* (median = 11.1 RPM, MAD = 12.2) was
222 markedly more abundant in MFMB-83 than any other sample (Figure 5B). These genera highlight
223 variability between microbiomes from a single food source and may provide insights into other
224 dissimilarities in these samples.

225

226 **2.5. Microbiome shifts in response to changes in food matrix composition**

227 We tested the hypothesis that the microbiome composition will shift in response to changes
228 in the food matrix and can be a unique signal to indicate contamination or adulteration. In 28 of
229 the 31 HPP samples, >99% of the matrix reads were determined in our related work³⁸ to originate
230 from poultry (*Gallus gallus*), which was the only ingredient expected based on ingredient
231 specifications. However, three samples had higher pork and beef content compared to all other
232 HPP samples: MFMB-04 (7.74% pork, 8.99% beef), MFMB-20 (0.53% pork, 1.00% beef), and
233 MFMB-38 (0.92% pork, 0.29% beef) compared to the highest pork (0.01%) and beef (0.00%)
234 content among the other 28 HPP samples (Supplementary Data by Haiminen *et al.*³⁸). The
235 microbiomes of these matrix contaminated samples also clustered into a separate sub-cluster
236 (Figure 2B). This demonstrated that a shift in the food matrix composition was associated with an
237 observable shift in the food microbiome.

238 MFMB-04 and MFMB-20 had the highest percentage of microbial reads compared to other
239 samples (Supplementary Figure S3). They also exhibited an increase in *Lactococcus*,
240 *Lactobacillus*, and *Streptococcus* relative abundances compared to other samples (Figure 5B), also
241 reflected at respective higher taxonomic levels above genus (Supplementary Figure S5).

242 There were 53 genera identified uniquely in MFMB-04 and/or MFMB-20, but not present
243 in any other sample. (MFMB-38 had a very low microbial load and contributed no uniquely
244 identified genera above the abundance threshold.) MFMB-04 contained 44 unique genera (Figure
245 4) with the most abundant being *Macrococcus* (35.8 RPM), *Psychrobacter* (23.8 RPM), and
246 *Brevibacterium* (18.1 RPM). Additionally, *Paenalcaligenes* was present only in MFMB-04 and
247 MFMB-20 with an RPM of 6.4 and 0.3, respectively, compared to a median RPM of 0.004 among
248 other samples. Notable differences in the matrix-contaminated samples' unique microbial
249 community membership compared to other samples may provide microbial indicators associated
250 with unanticipated pork or beef presence.

251 **2.6. Genus level identification of foodborne microbes**

252 We evaluated the ability of total RNA sequencing to identify genera of commonly known
253 foodborne pathogens within the microbiome. We focused on fourteen pathogen-containing genera
254 including *Aeromonas*, *Bacillus*, *Campylobacter*, *Clostridium*, *Corynebacterium*, *Cronobacter*,
255 *Escherichia*, *Helicobacter*, *Listeria*, *Salmonella*, *Shigella*, *Staphylococcus*, *Vibrio*, and *Yersinia*
256 that were found to be present in the HPP samples with varying relative abundances. Of these
257 genera, *Aeromonas*, *Bacillus*, *Campylobacter*, *Clostridium*, *Corynebacterium*, *Escherichia*,
258 *Salmonella*, and *Staphylococcus* were detected in every HPP with median abundance values
259 between 0.58–48.31 RPM (Figure 6A). This indicated that a baseline fraction of reads can be
260 attributed to foodborne microbes when using NGS. Of those genera appearing in all samples, there

261 was observed sample-to-sample variation in their abundance with some genera exhibiting longer
262 tails of high abundance, e.g. *Staphylococcus* and *Salmonella*, whereas others exhibit very low
263 abundance barely above the threshold of detection, e.g. *Bacillus* and *Yersinia* (Figure 6A). None
264 of the pathogen-containing genera were consistent with higher relative abundances due to
265 differences in food matrix composition. *Bacillus* and *Corynebacterium* exhibited slightly higher
266 relative abundances in sample MFMB-04 which contained 7.7% pork and 9.0% beef (Figure 6B).
267 Yet while MFMB-04 contained higher cumulative levels of these foodborne microbes, the next
268 highest sample was MFMB-93 which was not associated with altered matrix composition, and
269 both MFMB-04 and MFMB-93 contained higher levels of *Staphylococcus* (Figure 6B). Thus,
270 matrix composition alone did not explain variations of these pathogen-containing genera.

271 Interestingly, low to moderate levels of *Salmonella* were detected within all 31 HPP
272 microbiomes (Figure 6A). The presence of *Salmonella* in HPP is expected but the viability of
273 *Salmonella* is an important indicator of safety and quality. Thus, we further sought to delineate
274 *Salmonella* growth capability within these microbiomes by comparing culturability with multiple
275 established bioinformatic NGS methods for *Salmonella* relative abundances in the samples.

276 **2.7 Assessment of *Salmonella* culturability and total RNA sequencing**

277 Total RNA sequencing of food microbiomes has the potential to provide additional
278 sensitivity beyond standard culture-based food safety testing to confirm or reject the presence of
279 potentially pathogenic microbes. In all of the examined HPP samples, some portion of the
280 sequenced reads were classified as belonging to pathogen-containing genera (Figure 6); however,
281 the presence of RNA transcripts does not necessarily indicate current growth of the organism itself.
282 We further inspected one pathogen of interest, *Salmonella*, to determine the congruence between
283 sequencing-based and culturability results. Of the 31 samples examined with total RNA

284 sequencing, *Salmonella* culture testing was applied to 27 samples, of which four were culture-
285 positive. Surprisingly, *Salmonella* culture-positive samples were not among those with the highest
286 relative abundance of *Salmonella* from sequencing (Figure 7A). When ranking the samples by
287 decreasing *Salmonella* abundance, the culture-positive samples were not enriched for higher ranks
288 ($p=0.86$ from Wilcoxon rank sum test indicating that the distributions are not significantly
289 different, Table 2). To confirm that the microbiome analysis pipeline did not miss *Salmonella* reads
290 present, we completed two orthogonal analyses on the same data set used in the microbial
291 identification step. The reference genomes relevant to these additional analyses were publicly
292 available and closed high quality genomes available from the sources indicated below.

293 First, for a targeted analysis, we aligned the sequenced reads using a different tool, Bowtie
294 2,⁴⁸ to an augmented *Salmonella*-only reference database. This reference was comprised of the 264
295 *Salmonella* genomes extracted from NCBI RefSeq Complete (used in our previous microbial
296 identification step) as well as an additional 1,183 public *Salmonella* genomes which represent
297 global diversity within the genus.⁴⁹ The number of reads that aligned to the *Salmonella*-only
298 reference was on average 370-fold higher than identified as *Salmonella* by Kraken using the multi-
299 microbe NCBI RefSeq Complete. In this additional analysis, the culture-positive samples had
300 overall higher ranks compared to culture-negative samples ($p=0.06$, Table 2) indicating that
301 additional *Salmonella* genomic data in the reference significantly improved discriminatory
302 identification power. *Salmonella* culture-positive samples were still not the most abundant (Figure
303 7B), but with an enriched database, sequencing positioned all four culturable samples within the
304 top 10 ranking.

305 The second additional analysis examined alignment of the reads to a specific gene
306 required⁵⁰ for replication and protein production in actively dividing *Salmonella*— elongation
307 factor Tu (*ef-Tu*). This was done by aligning the reads to 4,846 gene sequences for *ef-Tu* extracted

308 for a larger corpus of *Salmonella* genomes from OMXWare.⁵¹ The relative abundances of this
309 transcript in culture-positive samples were still comparable to culture-negative samples (Figure
310 7C). Culture-positive samples did not exhibit higher ranks compared to culture-negative samples
311 ($p=0.56$, Table 2), indicating that *ef-Tu* relative abundance alone was not sufficient to improve the
312 lack of concordance in culturability vs sequencing. These two orthogonal analyses demonstrated
313 that results from carefully developed culture-based testing and those from current high-throughput
314 sequencing technologies, whether assessed at overall reads aligned or specific gene abundances,
315 were not conclusively in agreement when detecting active *Salmonella* in food samples (Figure 7
316 and Table 2). However, the use of a reference database enriched in whole genome sequences of
317 the specific organism of interested was found appropriate for food safety applications.

318 Since microbes compete for available resources within an environmental niche and
319 therefore impact one another,⁵² we investigated *Salmonella* culture results in conjunction with co-
320 occurrence patterns of other microbes in the total RNA sequencing data (Figure 8). Point-biserial
321 correlation coefficients (r_{pb}) were calculated between *Salmonella* culturability results (presence or
322 absence which were available for 27 of the 31 samples) and microbiome relative abundance. We
323 observed 31 genera that positively correlated and with *Salmonella* presence ($r_{pb} > 0.5$).
324 *Erysipelothrix*, *Lactobacillus*, *Anaerococcus*, *Brachyspira*, and *Jeotgalibaca* exhibited the largest
325 positive correlations. *Gyrovirus* was negatively correlated with *Salmonella* growth ($r_{pb} = -0.54$).
326 In three of the four *Salmonella*-positive samples (MFMB-04, MFMB-20, and MFMB-38), food
327 matrix contamination was also observed (Supplementary Data in Haiminen *et al.*³⁸). The
328 concurrency of *Salmonella* growth and matrix contamination was affirmed by the microbial co-
329 occurrence (specifically *Erysipelothrix*, *Brachyspira*, and *Gyrovirus*). This highlights the complex
330 dynamic and community co-dependency of food microbiomes, yet shows that multiple dimensions

331 of the data (microbiome composition, culture-based methods, and microbial load) will signal
332 anomalies from typical samples when there is an issue in the supply chain.

333

334 **3. DISCUSSION:**

335 Accurate and appropriate tests for detecting potential hazards in the food supply chain are key to
336 ensuring consumer safety and food quality. Monitoring and regular testing of raw ingredients can
337 reveal fluctuations within the supply chain that may be an indicator of an ingredient's quality or
338 of a potential hazard. Such quality is assessed by standardized tests for chemical and microbial
339 composition to meet legal requirements and specifications from government agencies throughout
340 the world. For raw materials or finished products to meet these bounds of safety and quality, their
341 composition must usually have a low microbiological load (except in fermented foods) and be
342 chemically identical in macro-components such as carbohydrate, protein, and fat. Methods in this
343 space must avoid false negative results which could endanger consumers, while also minimizing
344 false positives which could lead to unnecessary recalls and food loss.

345 Existing microbial detection technologies used in food safety today such as pulse field gel
346 electrophoresis (PFGE) and whole genome sequencing (WGS) require microbial isolation. This
347 provides biased outcomes as it removes microbes from their native environment where other biotic
348 members also subsist, and selects microbes by culturability alone. Amplicon sequencing, while a
349 low-cost alternative to metagenome or metatranscriptome sequencing for bacteria, also imparts
350 PCR amplification bias and reduces detection sensitivity due to reliance on a single gene (16S
351 ribosomal RNA).^{14,53,54} We therefore investigated the utility of total RNA sequencing of food
352 microbiomes and demonstrated that from this single test, we are able to yield several pertinent
353 results about food safety and quality.

354 For this evaluation, we developed a pipeline to characterize the microbiome of typical food
355 ingredient samples and to detect potentially hazardous outliers. Special considerations for food
356 samples were made as computational pipelines for human or other microbiome analyses are not
357 sufficient for applications in food safety without modification. In food, the eukaryotic matrix needs
358 to be confirmed, may be mixed, and, as we and others have shown, affects the identification
359 accuracy of microbes that are present.^{35,36} By filtering food matrix sequence data properly, we
360 avoid incorrect microbial identification and characterization of the microbiome³⁶ while also
361 increasing the computational efficiency for downstream processing. The addition of this filtering
362 step in the pipeline removed approximately 90% of false positive genera and provided results at
363 99.96% specificity when evaluating simulated mixtures of food matrix and microbes (Table 1).

364 Through the analysis of 31 high protein powder total RNA sequencing samples, we
365 demonstrated the pipeline's ability to characterize food microbiomes and indicate outliers. In this
366 sample collection, we identified a core catalog of 65 microbial genera found in all samples where
367 *Bacteroides*, *Clostridium*, and *Lactococcus* were the most abundant (Supplementary Table 4). We
368 also demonstrated that in these food microbiomes the overall diversity was 2-fold greater than the
369 core microbe set. Fluctuations in the microbiome can indicate important differences between
370 samples as observed here, as well as in the literature for grape berry⁶ and apple fruit microbiomes
371 (pertaining to organic versus conventional farming)⁷ or indicate inherent variability between
372 production batches or suppliers as observed here and during cheddar cheese manufacturing.⁸
373 Specifically, we observed a shift in the microbial composition (Figure 2B) and the microbial load
374 (Supplementary Figure S3) in high protein powder samples (derived from poultry meal) where
375 unexpected pork and beef were observed. Matrix-contaminated samples were marked by increased
376 relative abundances of specific microbes including *Lactococcus*, *Lactobacillus*, and *Streptococcus*
377 (Figure 5B). This work shows that the microbiome shifts with observed food matrix contamination

378 from sources with similar macronutrient content and thus, the microbiome alone is a likely signal
379 of compositional change in food.

380 Beyond shifts in the microbiome, we focused on a set of well-defined foodborne-pathogen
381 containing genera and explored their relative abundances observed from total RNA sequencing.
382 Of these genera, *Aeromonas*, *Bacillus*, *Campylobacter*, *Clostridium*, *Corynebacterium*,
383 *Escherichia*, *Salmonella*, and *Staphylococcus* were detected in every HPP sample. This highlights
384 that when using NGS there may be an observable baseline of sequences assigned to potentially
385 pathogenic microbes. For this ingredient type, this result lends a range of normalcy of relative
386 abundance generated by NGS. Further work is needed to establish a definitive and quantitative
387 range of typical variation in samples of a particular food source and the degree of anomaly for a
388 new sample or genus abundance. However, preliminary studies of this nature can inform the
389 development of guidelines when working with increasingly sensitive shotgun metagenomic or
390 metatranscriptomic analysis.

391 Furthermore, sequenced DNA or RNA alone does not imply microbial viability. Therefore,
392 we investigated the relatedness of culture-based tests and total RNA sequencing for the pathogenic
393 bacterium *Salmonella* in the high protein powder samples. As has been reported for human gut⁵⁵
394 and deep sea⁵⁶ microbiomes, we also did not detect a correlation between *Salmonella* read
395 abundance and culturability (Figure 7 and Table 2). Sequence reads matching *Salmonella*
396 references were observed for all samples (both culture-positive and culture-negative) as
397 determined by multiple analysis techniques: microbiome classification, alignment to *Salmonella*
398 genomes, and targeted growth gene analysis. When ranking the high protein powder samples based
399 on *Salmonella* abundance from whole genome alignments, the culture-positive samples were
400 enriched for higher ranks ($p = 0.06$). However, the culture-positive samples were still intermixed
401 in ranking with culture-negative samples. This indicated that there was no clear minimum

402 threshold of sequence data as evidence for culturability and that this analysis alone is not predictive
403 of pathogen growth. One possible reason for this is that the culture-positive variant of *Salmonella*
404 is missing from existing reference data sets. Potentially, *Salmonella* attained a nonculturable state
405 wherein it was detected by sequencing techniques yet remained nonculturable from the HPP
406 sources. Successful isolation of total RNA and DNA and gene expression analysis from
407 experimentally known nonculturable bacteria has been demonstrated by Ganesan *et al.* in multiple
408 studies in other genera.^{19,22} Physiological state should thus be taken under consideration when
409 benchmarking sequencing technologies in comparison with culture-based methods. Thus, total
410 RNA sequencing of food samples may identify shifts that standard food testing does not, but the
411 incongruity between sequencing read data and culture-based results highlights the need to perform
412 more benchmarking in food microbiome analysis for pathogen detection.

413 The characterization of HPP food microbiomes leveraged current accepted public reference
414 databases, yet it is known that these databases are still inadequate.^{1,2,11,57,58} Furthermore, when
415 considering congruence between *Salmonella* culturability and NGS read mapping techniques, the
416 genetic breadth and depth of multi-genome reference sequences is essential. For example, focusing
417 on *ef-Tu*, a known marker gene for *Salmonella* growth, was not sufficient to mirror viability of *in*
418 *vitro* culture tests. This highlights the limitations of single gene approaches for identification.
419 When the sequenced reads were examined in the context of an augmented reference collection of
420 *Salmonella* genomes, we observed improved ranking and read mapping rate for culture-positive
421 samples (yet we did not achieve complete concordance). This improvement underlined the
422 increased analytical robustness yielded from a multi-genome reference. We also recognize that the
423 read mapping rate may be exaggerated as reads from non-*Salmonella* genomes could map to
424 *Salmonella* in the absence of any other reference genomes. Overall for robust analysis and
425 applicability to food safety and quality, microbial references must be expanded to include more

426 genetically diverse representatives of pathogenic and spoilage organisms. Description of food
427 microbiomes will only improve as additional public sequence data is collected and leveraged.

428 In our sample collection, 2-4% (effectively 5 to 14 million) of reads remain unclassified. The
429 GC content distribution of unclassified reads matched microbial GC content distribution
430 (Supplementary Figure S4) suggesting that these reads may have been derived from microbes
431 missing from the current reference database that have not yet been isolated or sequenced. By
432 sequencing the microbiome, we sampled environmental niches in their native state in a culture-
433 independent manner and therefore collected data from diverse and potentially never-before seen
434 microbes. Tracking unclassified reads will also be essential for monitoring food microbiomes. The
435 inability to provide a name from existing references does not eliminate the possibility that the
436 sequence is from an unwanted microbe or indicates a hazard. In addition to tracking known
437 microbes, quantitative or qualitative shifts in the unclassified sequences might be used to detect
438 when a sample is different from its peers.

439 We demonstrated the potential utility of analyzing food microbiomes for food safety using raw
440 ingredients. This study resulted in the detection of shifts in the microbiome composition
441 corresponding to unexpected matrix contaminants. This signifies that the microbiome is likely an
442 important and effective hazard indicator in the food supply chain. While we have used total RNA
443 sequencing for detection of microbiome membership, the technology has future applicability for
444 detection of antimicrobial resistance, virulence, and biological function for multiple food sources,
445 and for other sample types. Notably, while this pipeline was developed for food monitoring, with
446 applicable modifications and identification of material-specific indicators, it can be applied to
447 other microbiomes including human and environmental.

448

449 **4. METHODS:**

450

451 **4.1 Sample Collection, Preparation, and Sequencing**

452 High protein powder (HPP, 2.5 kg) samples were each collected from a train car in Reno, NV,
453 USA between April 2015 and February 2016 in four batches from two suppliers and shipped to
454 the Weimer lab at the University of California, Davis (Davis, CA). Each HPP sample was
455 composed of five sub-samples from random locations within the train car prior to shipment.
456 Sample preparation, total RNA extraction and integrity confirmation, cDNA construction, and
457 library construction for these samples was previously described by Haiminen et al.³⁸

458 Sequencing was performed by BGI@UC Davis (Sacramento, CA) using Illumina HiSeq
459 4000 (San Diego, CA) with 150 paired end chemistry for each sample except the following: HiSeq
460 3000 with 150 paired end chemistry was used for MFMB-04 and MFMB-17. All total RNA
461 sequencing data are available via the 100K Pathogen Genome Project BioProject (PRJNA186441)
462 at NCBI (Supplementary Table 1).

463 For evaluation of total RNA sequencing for microbial classification in paired processing
464 steps, total RNA and total DNA were extracted from the same sample and denoted as MFMB-03
465 and MFMB-08, respectively. Total RNA was extracted and sequenced as described above. Total
466 DNA was extracted and sequenced as described previously.^{10,59-64} The Illumina HiSeq 2000 with
467 100 paired end chemistry was used for MFMB-03 and MFMB-08.

468

469 **4.2 Sequence Data Quality Control**

470 Illumina Universal adapters were removed and reads were trimmed using Trim Galore⁶⁵
471 with a minimum read length parameter 50 bp. The resulting reads were filtered using Kraken³⁷, as
472 described below in Section 4.3, with a custom database built from the PhiX genome (NCBI
473 Reference Sequence: NC_001422.1). Removal of PhiX content is suggested as it is a common

474 contaminant in Illumina sequencing data.⁶⁶ Trimmed non-PhiX reads were used in subsequent
475 matrix filtering and microbial identification steps.

476

477 **4.3 Matrix Filtering Process and Validation**

478 Kraken³⁷ with a *k*-mer size of 31 bp (optimal size described in the Kraken reference
479 publication) was used to identify and remove reads that matched a pre-determined list of 31
480 common food matrix and potential contaminant eukaryotic genomes (Supplementary Table 2).
481 These food matrix organisms were chosen based on preliminary eukaryotic read alignment
482 experiments of the HPP samples as well as high-volume food components in the supply chain. Due
483 to the large size of eukaryotic genomes in the custom Kraken³⁷ database, a random *k*-mer reduction
484 was applied to reduce the size of the database by 58% using kraken-build with option --max-db-
485 size, in order to fit the database in 188 GB for in-memory processing. A conservative Kraken score
486 threshold of 0.1 was applied to avoid filtering microbial reads. The matrix filtering database
487 includes low complexity and repeat regions of eukaryotic genomes to capture all possible matrix
488 reads. This filtering database with the score threshold was also used in the matrix filtering *in silico*
489 testing as described below.

490 Matrix filtering was validated by constructing synthetic paired end reads (150 bp) using
491 DWGSIM⁶⁷ with mutations from reference sequences using the following parameters: base error
492 rate (*e*) = 0.005, outer distance between the two ends of a read pair (*d*) = 500, rate of mutations (*r*)
493 = 0.001, fraction of indels (*R*) = 0.15, probability an indel is extended (*X*) = 0.3. Reference
494 sequences are detailed in Supplementary Table 3. We constructed two *in silico* mixtures of
495 sequencing reads by randomly sampling reads from eukaryotic reference genomes. Simulated
496 Food Mixture 1 was comprised of nine species with the following number of reads per genome:
497 2M cattle, 2M salmon, 1M goat, 1M lamb, 1M tilapia (transcriptome), 962K chicken

498 (transcriptome), 10K duck, 1K horse, and 1K rat totaling 7.974M matrix reads. Simulated Food
499 Mixture 2 contained 5M soybean, 4M rice, 3M potato, 2M corn, 200K rat, and 10K drain fly reads,
500 totaling 14.210M matrix reads. Both simulated food mixtures included 1,000 microbial sequence
501 reads generated from 15 different microbial species for a total of 15K sequence reads
502 (Supplementary Table 3).

503

504 **4.4 Microbial Identification**

505 Remaining reads after quality control and matrix filtering were classified using Kraken³⁷
506 against a microbial database with a k -mer size of 31 bp to determine the microbial composition
507 within each sample. NCBI RefSeq Complete⁶⁸ genomes were obtained for bacterial, archaeal,
508 viral, and eukaryotic microorganisms (~7,800 genomes retrieved April 2017). Low complexity
509 regions of the genomes were masked using Dustmasker⁶⁹ with default parameters. A threshold of
510 0.05 was applied to the Kraken score in an effort to maximize the F-score of the result (as
511 demonstrated in Kraken's operating manual.⁷⁰ Taxa-specific sequence reads were used to calculate
512 a relative abundance in reads per million (RPM; Equation 1) where R_T represents the reads
513 classified per microbial entity (e.g. the genus *Salmonella*) and R_Q represents the number of
514 sequenced reads remaining after quality control (trimming and PhiX removal) for an individual
515 sample, including any reads classified as eukaryotic:

516

$$517 \quad RPM = \frac{R_T}{R_Q} \times 1,000,000 \quad \text{Equation 1}$$

518

519 This value provides a relative abundance of the microbial entity of interest and was used in
520 comparisons of taxa among samples. Genera with a conservative threshold of $RPM > 0.1$ were

521 defined as present, as previously applied by others in the contexts of human infectious disease and
522 gut microbiome studies.^{33,34} Pearson correlation of resulting microbial genus counts was
523 computed.

524

525 **4.5 Community Ecology Analysis**

526 Rarefaction analysis at multiple subsampled read depths R_D was performed by multiplying
527 the microbial genus read counts with R_D/R_Q and rounding the results down to the nearest integer
528 to represent observed read counts. Here R_Q is the total number of reads in the sample after quality
529 control (including microbial, matrix, and unclassified reads). Resulting α -diversity at read depth
530 R_D was computed as the number of genera with resulting RPM > 0.1 and plotted at five million
531 read intervals: $R_D = 5M, 10M, 15M, \dots, R_Q$. If, due to random sampling and rounding effects, the
532 computed α -diversity was lower than the diversity computed at any previous depth, the previous
533 higher α -diversity was used for plotting. The median elbow was calculated as previously
534 described⁴⁵ using the R package `kneed`.⁴⁵

535 In compositional data analysis,³¹ non-zero values are required when computing β -diversity
536 based on Aitchison distance.⁴⁶ Therefore, reads counts assigned to each genus were pseudo-
537 counted by adding one in advance of computation of RPM (Eq. 1) prior to calculating the Aitchison
538 distance for the microbial table. β -diversity was calculated using the R package `robCompositions`⁷¹
539 and hierarchical clustering was performed using base R function `hclust` using the “ward.D2”
540 method as recommended for compositional data analysis.³¹

541

542 **4.6 Unclassified Read Analysis**

543 The GC percent distributions of matrix (from matrix filtering), microbial, and remaining
544 unclassified reads per sample were computed using FastQC⁷² and collated across samples with
545 MultiQC.⁷³

546

547 **4.7 Analysis of *Salmonella* Culturability**

548 Growth of *Salmonella* was determined using a real-time quantitative PCR method for the
549 confirmation of *Salmonella* isolates for presumptive generic identification of foodborne
550 *Salmonella*. Testing was performed fully in concordance with the Bacteriological Analytical
551 Manual (BAM) for *Salmonella*^{74,75} for this approach that is also AOAC-approved. All samples
552 with positive results for *Salmonella* were classified as containing actively growing *Salmonella*. To
553 compare culture results with those from total RNA sequencing, *Salmonella* RPM values were
554 parsed from the genus-level microbe table (described in Section 4.4).

555 Two additional approaches were employed to examine *Salmonella* read mapping with a
556 more sensitive tool and broader reference databases. Quality controlled matrix-filtered reads were
557 aligned using Bowtie2⁴⁸ with very-sensitive-local-mode to 1. an expanded collection of whole
558 *Salmonella* genomes and 2. to a curated growth gene reference for elongation factor Tu (*ef-Tu*).
559 For results from both complete genome and *ef-Tu* gene alignments, the relative abundance (RPM)
560 was computed as shown in Equation 1.

561 For whole genome alignments, a reference was constructed from 1,183 recently published
562 *Salmonella* genomes⁴⁹ in addition to the 264 *Salmonella* genomes extracted from the
563 aforementioned NCBI RefSeq Complete collection (see Methods Section 4.4).

564 To construct a curated growth gene (*ef-Tu*) reference, gene sequences annotated in
565 *Salmonella* genomes as “elongation factor Tu”, “EF-Tu” or “eftu” (case insensitive) were retrieved
566 from OMXWare⁵¹ using its Python package. This query yielded 4,846 unique gene sequences from

567 a total of 36,242 *Salmonella* genomes which were assembled or retrieved from the NCBI Sequence
568 Read Archive or RefSeq Complete Sequences as previously described.⁵¹ The retrieved *ef-Tu* gene
569 sequences were subsequently used to build a custom Bowtie2⁴⁸ reference. Read alignment was
570 completed with very-sensitive-local-mode.

571 The read counts for each sample were ranked and Wilcoxon rank sum test was computed
572 between the rank vectors of 4 *Salmonella*-positive and 23 *Salmonella*-negative samples. The 4
573 samples with unknown *Salmonella* status were excluded from the rankings.

574 Point-biserial correlation coefficients (r_{pb}) were calculated between *Salmonella* growth
575 indicated by culture results (+1 and -1 for presence and absence, respectively) and observed
576 relative abundance from total RNA sequencing results using the R package ltm.⁷⁶ The point-
577 biserial correlation is a special case of the Pearson correlation that is better suited for a binary
578 variable e.g. when *Salmonella* is reported as present or absent (a sample's *Salmonella* status).

579

580 **Data Availability:**

581 All high protein powder (HPP) poultry meal sequences are available through the 100K
582 Pathogen Genome Project (PRJNA186441) in the NCBI BioProject (see Supplementary Table 1
583 for a complete list of accession numbers).

584

585 **Code Availability:**

586 The pipeline and microbial or matrix references were constructed from publicly available
587 tools and reference sequences as described in the Methods section. Automated usability of this
588 pipeline is available through membership in the Consortium for Sequencing the Food Supply
589 Chain.

590

591 **Acknowledgements:**

592 We'd like to acknowledge the IBM Research OMXWare team for their data management
593 support and availability for the retrieval and processing of microbial genomes. This research
594 project was financially supported by the Consortium for Sequencing the Food Supply Chain.
595 Funding for the total RNA sequencing of high protein powder factory ingredients was provided by
596 Mars, Incorporated to B.C.W. with specific interest in metagenomics of the food microbiome.

597

598 **Contributions:**

599 KLB and NH conceived of the experimental design, developed the approach, completed
600 and oversaw the experiments, performed analyses, and wrote the paper; DC, SE, MK, BK, MD,
601 RP, HK, ES developed the approach, analyzed data, and revised the manuscript; BCH completed
602 nucleic acid extraction method development and sequencing library construction, and contributed
603 to data analysis and writing; NK coordinated sample collection and processing, nucleic acid
604 extraction and contributed to writing; RB and PM conceived of the experimental design, developed
605 the approach, and reviewed the paper; BG contributed to the experimental design, developed the
606 approach, and wrote the paper; GD, CHM, SP, AQ participated to the conception of the
607 experimental design and to the review of the manuscript; LP conceived of the experiment,
608 contributed to the data analysis, and wrote the paper; JHK conceived of the experiment, developed
609 the approach, and wrote the paper; BCW conceived of the experimental design, developed the
610 approach, oversaw the experiments, performed analyses, and wrote the paper

611

612 **Competing Interests:**

613 The authors were employed by private or academic organizations as described in the author
614 affiliations at the time this work was completed. IBM Corporation, Mars Incorporated, and Bio-

615 Rad Laboratories are members of the Consortium for Sequencing the Food Supply Chain. The
616 authors declare no other competing interests

617

618 **Supplementary information is available at npj Science of Food's website**

619

620 REFERENCES:

- 621 1. Kovac, J., Bakker, H. den, Carroll, L. M. & Wiedmann, M. Precision food safety: A
622 systems approach to food safety facilitated by genomics tools. *TrAC Trends Anal. Chem.*
623 (2017). doi:10.1016/j.trac.2017.06.001
- 624 2. Weimer, B. C. *et al.* Defining the food microbiome for authentication, safety, and process
625 management. *IBM J. Res. Dev.* **60**, 1 (2016).
- 626 3. Walsh, A. M. *et al.* Microbial Succession and Flavor Production in the Fermented Dairy
627 Beverage Kefir. *mSystems* **1**, (2016).
- 628 4. Walsh, A. M. *et al.* Species classifier choice is a key consideration when analysing low-
629 complexity food microbiome data. *Microbiome* **6**, 50 (2018).
- 630 5. Duru, I. C. *et al.* Metagenomic and metatranscriptomic analysis of the microbial
631 community in Swiss-type Maasdam cheese during ripening. *Int. J. Food Microbiol.* **281**,
632 10–22 (2018).
- 633 6. Martins, G. *et al.* Grape berry bacterial microbiota: Impact of the ripening process and the
634 farming system. *Int. J. Food Microbiol.* **158**, 93–100 (2012).
- 635 7. Abdelfattah, A., Wisniewski, M., Droby, S. & Schena, L. Spatial and compositional
636 variation in the fungal communities of organic and conventionally grown apple fruit at the
637 consumer point-of-purchase. *Hortic. Res.* **3**, 16047 (2016).
- 638 8. Williams, A. G., Choi, S.-C. & Banks, J. M. Variability of the species and strain
639 phenotype composition of the non-starter lactic acid bacterial population of cheddar
640 cheese manufactured in a commercial creamery. *Food Res. Int.* **35**, 483–493 (2002).
- 641 9. Weimer, B. C. 100K Pathogen Genome Project. *Genome Announc.* **5**, e00594-17 (2017).
- 642 10. Emond-Rheault, J.-G. *et al.* A Syst-OMICS Approach to Ensuring Food Safety and
643 Reducing the Economic Burden of Salmonellosis. *Front. Microbiol.* **8**, 996 (2017).
- 644 11. Kaufman, J. H. *et al.* Insular microbiogeography. (2017).
- 645 12. Bashiardes, S., Zilberman-Schapira, G. & Elinav, E. Use of Metatranscriptomics in
646 Microbiome Research. *Bioinform. Biol. Insights* **10**, 19–25 (2016).
- 647 13. McGrath, K. C. *et al.* Isolation and analysis of mRNA from environmental microbial
648 communities. *J. Microbiol. Methods* **75**, 172–176 (2008).
- 649 14. Cottier, F. *et al.* Advantages of meta-total RNA sequencing (MeTRS) over shotgun
650 metagenomics and amplicon-based sequencing in the profiling of complex microbial
651 communities. *npj Biofilms Microbiomes* **4**, 2 (2018).
- 652 15. Macklaim, J. M. *et al.* Comparative meta-RNA-seq of the vaginal microbiota and
653 differential expression by *Lactobacillus iners* in health and dysbiosis. *Microbiome* **1**, 12
654 (2013).
- 655 16. Haiminen, N. *et al.* Food authentication from shotgun sequencing reads with an
656 application on high protein powders. *npj Sci. Food* **3**, (2019).

- 657 17. Lakin, S. M. *et al.* MEGARes: an antimicrobial resistance database for high throughput
658 sequencing.
- 659 18. Noyes, N. R. *et al.* Resistome diversity in cattle and the environment decreases during
660 beef production. *Elife* **5**, e13195 (2016).
- 661 19. Ganesan, B., Dobrowolski, P. & Weimer, B. C. Identification of the Leucine-to-2-
662 Methylbutyric Acid Catabolic Pathway of *Lactococcus lactis*. *Appl. Environ. Microbiol.*
663 **72**, 4264–4273 (2006).
- 664 20. Ganesan, B., Seefeldt, K., Koka, R. C., Dias, B. & Weimer, B. C. Monocarboxylic acid
665 production by lactococci and lactobacilli. *Int. Dairy J.* **14**, 237–246 (2004).
- 666 21. Ganesan, B., Seefeldt, K. & Weimer, B. C. Fatty Acid Production from Amino Acids
667 and -Keto Acids by *Brevibacterium linens* BL2. *Appl. Environ. Microbiol.* **70**, 6385–6393
668 (2004).
- 669 22. Ganesan, B., Stuart, M. R. & Weimer, B. C. Carbohydrate Starvation Causes a
670 Metabolically Active but Nonculturable State in *Lactococcus lactis*. *Appl. Environ.*
671 *Microbiol.* **73**, 2498–2512 (2007).
- 672 23. Ganesan, B. *et al.* Probiotic bacteria survive in Cheddar cheese and modify populations of
673 other lactic acid bacteria. *J. Appl. Microbiol.* **116**, 1642–1656 (2014).
- 674 24. Ganesan, B. & Weimer, B. C. *Cheese : chemistry, physics, and microbiology*. (Elsevier,
675 2004).
- 676 25. Sheflin, A. M., Melby, C. L., Carbonero, F. & Weir, T. L. Linking dietary patterns with
677 gut microbial composition and function. *Gut Microbes* **8**, (2017).
- 678 26. McDonald, D. *et al.* American Gut: an Open Platform for Citizen Science Microbiome
679 Research. *mSystems* **3**, e00031-18 (2018).
- 680 27. Clemente, J. C., Ursell, L. K., Parfrey, L. W. & Knight, R. The Impact of the Gut
681 Microbiota on Human Health: An Integrative View. *Cell* **148**, 1258–1270 (2012).
- 682 28. Richards, J. L., Yap, Y. A., McLeod, K. H., Mackay, C. R. & Mariño, E. Dietary
683 metabolites and the gut microbiota: an alternative approach to control inflammatory and
684 autoimmune diseases. *Clin Trans Immunol* **5**, e82 (2016).
- 685 29. Yang, X. *et al.* Use of Metagenomic Shotgun Sequencing Technology To Detect
686 Foodborne Pathogens within the Microbiome of the Beef Production Chain. *Appl Env.*
687 *Microbiol* **82**, 2433–2443 (2016).
- 688 30. Hofacre, C. L. *et al.* Characterization of antibiotic-resistant bacteria in rendered animal
689 products. *Avian Dis.* **45**, 953–61
- 690 31. Gloor, G. B., Macklaim, J. M., Pawlowsky-Glahn, V. & Egozcue, J. J. Microbiome
691 Datasets Are Compositional: And This Is Not Optional. *Front. Microbiol.* **8**, 2224 (2017).
- 692 32. Gloor, G. B. & Reid, G. Compositional analysis: a valid approach to analyze microbiome
693 high-throughput sequencing data. *Can. J. Microbiol.* **62**, 692–703 (2016).
- 694 33. Langelier, C. *et al.* Integrating host response and unbiased microbe detection for lower
695 respiratory tract infection diagnosis in critically ill adults. *Proc. Natl. Acad. Sci. U. S. A.*
696 **115**, E12353–E12362 (2018).
- 697 34. Ilott, N. E. *et al.* Defining the microbial transcriptional response to colitis through
698 integrated host and microbiome profiling. *ISME J.* **10**, 2389–2404 (2016).
- 699 35. Ripp, F. *et al.* All-Food-Seq (AFS): a quantifiable screen for species in biological samples
700 by deep DNA sequencing. *BMC Genomics* **15**, 639 (2014).
- 701 36. Lee, A. Y., Lee, C. S. & Gelder, R. N. Van. Scalable metagenomics alignment research
702 tool (SMART): a scalable, rapid, and complete search heuristic for the classification of
703 metagenomic sequences from complex sequence populations. *BMC Bioinformatics* **17**,
704 292 (2016).

- 705 37. Wood, D. E. & Salzberg, S. L. Kraken: ultrafast metagenomic sequence classification
706 using exact alignments. *Genome Biol.* **15**, R46 (2014).
- 707 38. Haiminen, N. *et al.* Food authentication from shotgun sequencing reads with an
708 application on high protein powders. *npj Sci. Food* in press (2019).
- 709 39. Hug, L. A. *et al.* A new view of the tree of life. *Nat. Microbiol.* **1**, 16048 (2016).
- 710 40. Wu, D. *et al.* A phylogeny-driven genomic encyclopaedia of Bacteria and Archaea.
711 *Nature* **462**, 1056–1060 (2009).
- 712 41. Kyrpides, N. C. *et al.* Genomic Encyclopedia of Bacteria and Archaea: Sequencing a
713 Myriad of Type Strains. *PLoS Biol.* **12**, e1001920 (2014).
- 714 42. Kyrpides, N. C., Eloe-Fadrosh, E. A. & Ivanova, N. N. Microbiome Data Science:
715 Understanding Our Microbial Planet. *Trends Microbiol.* **24**, 425–427 (2016).
- 716 43. Thompson, L. R. *et al.* A communal catalogue reveals Earth’s multiscale microbial
717 diversity. *Nature* **551**, 457 (2017).
- 718 44. Nayfach, S. & Pollard, K. S. Toward Accurate and Quantitative Comparative
719 Metagenomics. *Cell* **166**, 1103–1116 (2016).
- 720 45. Satopää, V., Albrecht, J., Irwin, D. & Raghavan, B. *Finding a ‘Kneedle’ in a Haystack:*
721 *Detecting Knee Points in System Behavior.*
- 722 46. Aitchison, J., Barceló-Vidal, C., Martín-Fernández, J. A. & Pawłowsky-Glahn, V.
723 Logratio Analysis and Compositional Distance. *Math. Geol.* **32**, 271–275 (2000).
- 724 47. Di Palma, M. A. & Gallo, M. A co-median approach to detect compositional outliers. *J.*
725 *Appl. Stat.* **43**, 2348–2362 (2016).
- 726 48. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods*
727 **9**, 357–9 (2012).
- 728 49. Kong, N. *et al.* Draft Genome Sequences of 1,183 Salmonella Strains from the 100K
729 Pathogen Genome Project. *Genome Announc.* **5**, (2017).
- 730 50. Tubulekas, I. & Hughes, D. A Single Amino Acid Substitution in Elongation Factor Tu
731 Disrupts Interaction between the Ternary Complex and the Ribosome. *J. Bacteriol.* 240–
732 250 (1993).
- 733 51. Seabolt, E. E. *et al.* OMXWare, A Cloud-Based Platform for Studying Microbial Life at
734 Scale. *arXiv* **1911.02095**, (2019).
- 735 52. Zelezniak, A. *et al.* Metabolic dependencies drive species co-occurrence in diverse
736 microbial communities. *Proc. Natl. Acad. Sci.* **112**, 6449–6454 (2015).
- 737 53. Jones, M. B. *et al.* Library preparation methodology can influence genomic and functional
738 predictions in human microbiome research. *Proc Natl Acad Sci U S A* (2015).
739 doi:10.1073/pnas.1519288112
- 740 54. Pollock, J., Glendinning, L., Wisedchanwet, T. & Watson, M. The madness of
741 microbiome: Attempting to find consensus ‘best practice’ for 16S microbiome studies.
742 *Appl. Environ. Microbiol.* AEM.02627-17 (2018). doi:10.1128/AEM.02627-17
- 743 55. Browne, H. P. *et al.* Culturing of ‘unculturable’ human microbiota reveals novel taxa and
744 extensive sporulation. *Nature* **533**, 543–546 (2016).
- 745 56. Eilers, H., Pernthaler, J., Glöckner, F. O. & Amann, R. Culturability and In situ abundance
746 of pelagic bacteria from the North Sea. *Appl. Environ. Microbiol.* **66**, 3044–51 (2000).
- 747 57. Hinchliff, C. E. *et al.* Synthesis of phylogeny and taxonomy into a comprehensive tree of
748 life. *Proc Natl Acad Sci U S A* **112**, 12764–12769 (2015).
- 749 58. Knight, R. *et al.* Best practices for analysing microbiomes. *Nat. Rev. Microbiol.* **16**, 410–
750 422 (2018).
- 751 59. Weis, A. M. *et al.* Genomic Comparison of *Campylobacter* spp. and Their Potential for
752 Zoonotic Transmission between Birds, Primates, and Livestock. *Appl. Environ. Microbiol.*

- 753 **82**, 7165 LP – 7175 (2016).
- 754 60. Miller, B. *et al.* A novel, single-tube enzymatic fragmentation and library construction
755 method enables fast turnaround times and improved data quality for microbial whole-
756 genome sequencing. *Kapa Biosyst. Appl. Note* 1–8 (2015).
757 doi:10.13140/RG.2.1.4534.3440
- 758 61. Lüdeke, C. H. M., Kong, N., Weimer, B. C., Fischer, M. & Jones, J. L. Complete genome
759 sequences of a clinical isolate and an environmental isolate of *Vibrio parahaemolyticus*.
760 *Genome Announc.* **3**, e00216-15 (2015).
- 761 62. Jeannotte, R. *et al.* High-Throughput Analysis of Foodborne Bacterial Genomic DNA
762 Using Agilent 2200 TapeStation and Genomic DNA ScreenTape System. *Agil. Appl. Note*
763 1–8 (2015). doi:doi:10.6084/m9.figshare.1372504
- 764 63. Arabyan, N. *et al.* Salmonella Degrades the Host Glycocalyx Leading to Altered Infection
765 and Glycan Remodeling. *Sci. Rep.* **6**, 1–11 (2016).
- 766 64. Kong, N. *et al.* Draft Genome Sequences of 1,183 *Salmonella* Strains from the 100K
767 Pathogen Genome Project. *Genome Announc.* **5**, e00518-17 (2017).
- 768 65. Krueger, F. TrimGalore. (2018).
- 769 66. Mukherjee, S., Huntemann, M., Ivanova, N., Kyrpides, N. C. & Pati, A. Large-scale
770 contamination of microbial isolate genomes by Illumina PhiX control. *Stand. Genomic*
771 *Sci.* **10**, 18 (2015).
- 772 67. Homer, N. DWGSIM. (2011).
- 773 68. O’Leary, N. A. *et al.* Reference sequence (RefSeq) database at NCBI: current status,
774 taxonomic expansion, and functional annotation. *Nucleic Acids Res.* **44**, D733–D745
775 (2016).
- 776 69. Morgulis, A., Gertz, E. M., Schäffer, A. A. & Agarwala, R. A fast and symmetric DUST
777 implementation to mask low-complexity DNA sequences. *J. Comput. Biol.* **13**, 1028–1040
778 (2006).
- 779 70. Wood, D. Kraken’s operating manual.
- 780 71. Templ, M., Hron, K. & Filzmoser, P. robCompositions: An R-package for Robust
781 Statistical Analysis of Compositional Data. in *Compositional Data Analysis* 341–355
782 (John Wiley & Sons, Ltd, 2011). doi:10.1002/9781119976462.ch25
- 783 72. Andrews, S. FastQC.
- 784 73. Ewels, P., Magnusson, M., Lundin, S. & Käller, M. MultiQC: summarize analysis results
785 for multiple tools and samples in a single report. *Bioinformatics* **32**, 3047–3048 (2016).
- 786 74. Andrews, W. H., Wang, H., Jacobson, A. & Hammack, T. Bacteriological Analytical
787 Manual (BAM) Chapter 5: Salmonella. in *Bacteriological Analytical Manual* (U.S. Food
788 and Drug Administration, 2018).
- 789 75. Grim, C. J. *et al.* High-Resolution Microbiome Profiling for Detection and Tracking of
790 *Salmonella enterica*. *Front. Microbiol.* **8**, 1587 (2017).
- 791 76. Rizopoulos, D. **Itm** : An R Package for Latent Variable Modeling and Item Response
792 Theory Analyses. *J. Stat. Softw.* **17**, 1–25 (2006).

793
794
795
796
797
798
799
800

801
802
803 **FIGURE and TABLE LEGENDS: (corresponding to their order at end of merged document)**
804
805 **Figure 1:** Pipeline description of bioinformatic steps applied to high protein powder
806 metatranscriptome samples. Black arrows indicate data flow and blue boxes describe outputs
807 from the pipeline.
808
809 **Table 1:** Accuracy of microbial identification using *in silico* constructed Simulated Food
810 Mixtures with expected food matrix and microbial sequences.
811
812 **Figure 2A:** Alpha diversity (number of genera) for all (n = 31) high protein powder
813 metatranscriptomes is compared to total number of sequenced reads for a range of *in silico*
814 subsampled sequencing depths. The dashed vertical line indicates the median elbow (at approx.
815 67 million reads).
816
817 **Figure 2B:** Hierarchical clustering of Aitchison distance values of poultry meal samples based
818 on microbial composition. Samples were received from Supplier A (blue and red) and Supplier B
819 (green). Matrix-contaminated samples are additionally marked in red.
820
821 **Figure 3A:** Phylogram of the 65 microbial genera present in all samples with RPM > 0.1
822
823 **Figure 3B:** Phylogram of all microbes observed in *any* sample. Log of the median RPM value
824 across samples is indicated. Grey indicating a median RPM value of 0.
825

826

827 **Figure 4:** Heatmap (\log_{10} -scale) of high protein powder microbial composition and relative
828 abundance (RPM) where absence ($\text{RPM} < 0.1$) is indicated in grey. Genera are ordered by
829 summed abundance across samples. Samples were received from Supplier A (blue) and Supplier
830 B (green). Red stars indicate matrix-contaminated samples (from Supplier A).

831

832 **Figure 5A:** All identified microbial general are plotted with median value and median absolute
833 deviation (MAD) of RPM abundance. Genera with $\text{MAD} > 5$ are labeled with the genus name.

834

835 **Figure 5B:** Heatmap (\log_{10} -scale) of ten microbial genera with the largest median absolute
836 deviation (MAD) across samples. Genera are ordered by decreasing MAD from top to bottom.
837 Samples were received from Supplier A (blue) and Supplier B (green). Red stars indicate matrix
838 contaminated samples (from Supplier A).

839

840 **Figure 6A:** Relative abundance of microbes with high relevance to food safety and quality from
841 high protein powder total RNA sequenced microbiomes. Width of violin plot indicates density of
842 samples with relative abundance at that value. Observation threshold of $\text{RPM} = 0.1$ is indicated
843 with the horizontal black line.

844

845 **Figure 6B:** Foodborne microbe relative abundances are shown across samples of high protein
846 powder total RNA sequenced samples.

847

848 **Figure 7:** *Salmonella* culturability status and high-throughput sequencing read abundance
849 (RPM) from *k*-mer classification to NCBI Microbial RefSeq Complete (**A**), from alignments to

850 1,447 *Salmonella* genomes (**B**), and from alignments to 4,846 EF-Tu gene sequences (**C**).
851 *Salmonella* presence (red) indicates culture-positive result, absence (green) indicates culture-
852 negative result, and no record (black) indicates samples for which no culture test was completed.

853
854 **Table 2:** The ranks for *Salmonella*-positive samples and the associated p-values from Wilcoxon
855 rank sum test are shown for high-throughput sequencing read abundance (RPM) for multiple
856 analyses: *k*-mer classification to NCBI Microbial RefSeq Complete (left), alignments to 1,447
857 *Salmonella* genomes (middle), and alignments to 4,846 *ef-Tu* gene sequences (right). The
858 corresponding *Salmonella* relative abundances are shown in Figure 7A–C.

859
860 **Figure 8:** *Salmonella* status correlations with genus relative abundances. Only those genera with
861 absolute value of the correlation coefficient > 0.5 are shown. Positive and negative correlations
862 are indicated in grey and blue, respectively.

863

864

865 **SUPPLEMENTAL INFORMATION:**

866 Supplemental Figures (pdf): Supplemental Figures S1–S5

867 Supplemental Table 1 (.xlsx) - Sample Descriptions

868 Supplemental Table 2 (.xlsx) - Matrix Filtering Genomes

869 Supplemental Table 3 (.xlsx) - Simulated Food Mixtures

870 Supplemental Table 4 (.xlsx) - Microbial Genera

RNA
sequencing
Reads
(FASTQ)

**Sequence
Quality Control**
TrimGalore
PhiX Filtering



Matrix Filtering
Classification with
Common Food
Ingredient
Database



**Microbial
Identification**
Classification with
microbial RefSeq
Complete

**Microbial
Quantification**
Genus-level relative
abundance (Eq. 1)

**Food Matrix
Reads**

**Unclassified
Reads**

**Relative
Abundance
(RPM)**

**Comparative
Statistics**
Intrasample
Intersample
Relative change

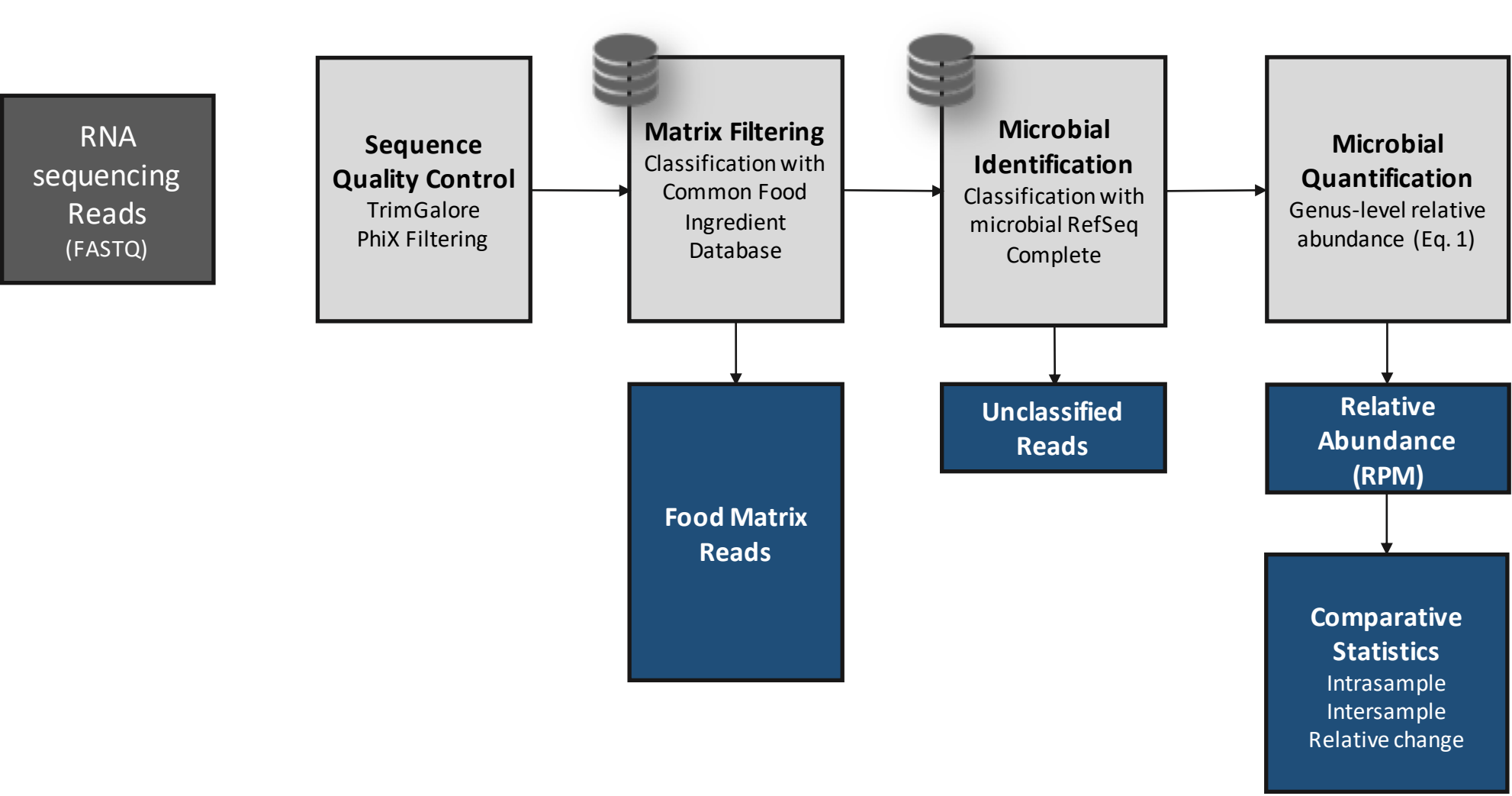
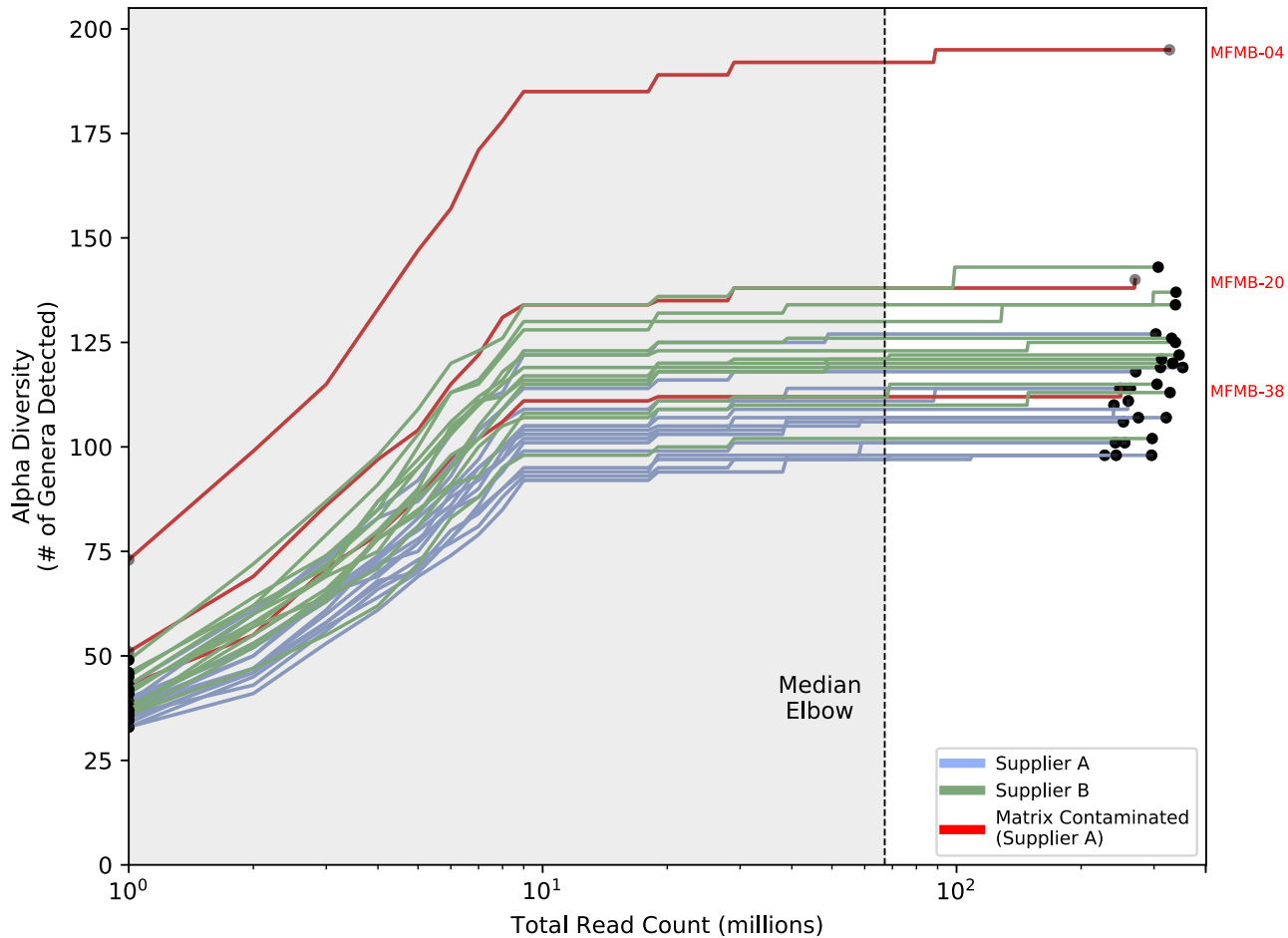
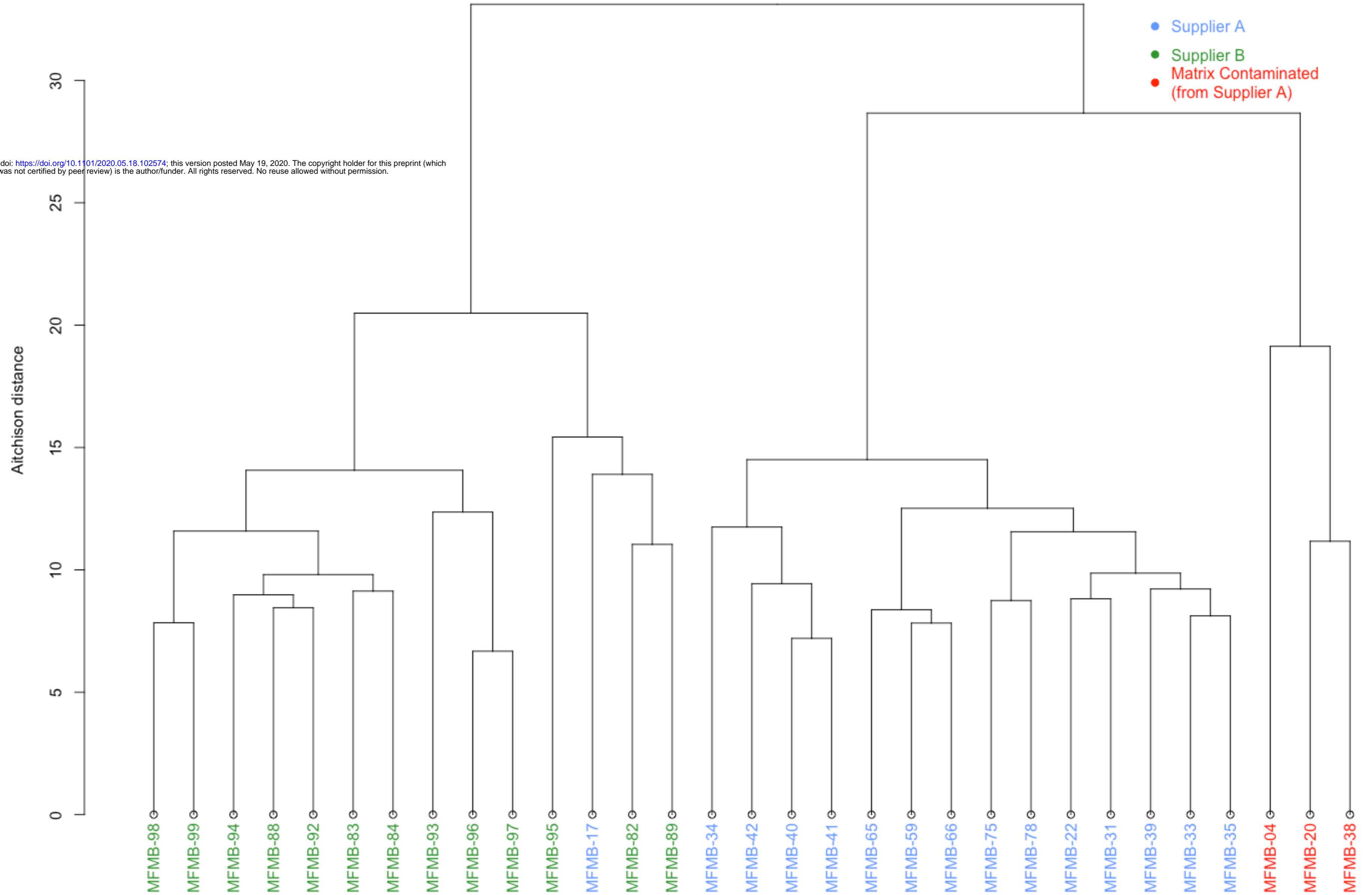


Table 1: Microbial Identification Accuracy from Simulated Food Microbiome Mixtures

	Simulated Mixture 1					Simulated Mixture 2			
	With Matrix Filtering		No Matrix Filtering			With Matrix Filtering		No Matrix Filtering	
	# GENERA	GENUS READS	# GENERA	GENUS READS		# GENERA	GENUS READS	# GENERA	GENUS READS
Bacteria in Simulated Mixture (Expected Content)	14	15,000	14	15,000		14	15,000	14	15,000
Observed Microbial Content									
Bacteria	18	13,517	34	13,700		15	13,551	33	13,999
Viruses	0	0	9	563		0	0	4	328
Archaea	0	0	1	1		0	0	1	3
Eukaryota	0	0	4	104		0	0	4	799
Total Observed	18	13,517	48	14,368		15	13,551	42	15,129
True Positives (as a % of total observed)	14 (78%)	13,511 (99.96%)	14 (29%)	13,571 (94.45%)		14 (93%)	13,548 (99.98%)	14 (33%)	13,623 (90.05%)
False Positives (as a % of total observed)	4 (22%)	6 (0.04%)	34 (71%)	797 (5.55%)		1 (7%)	3 (0.02%)	28 (67%)	1,506 (9.95%)
False Positives Removed with Matrix Filtering (as a % of false positives without filtering)	30 (88.2%)	791 (99.2%)				27 (96.4%)	1,503 (99.8%)		

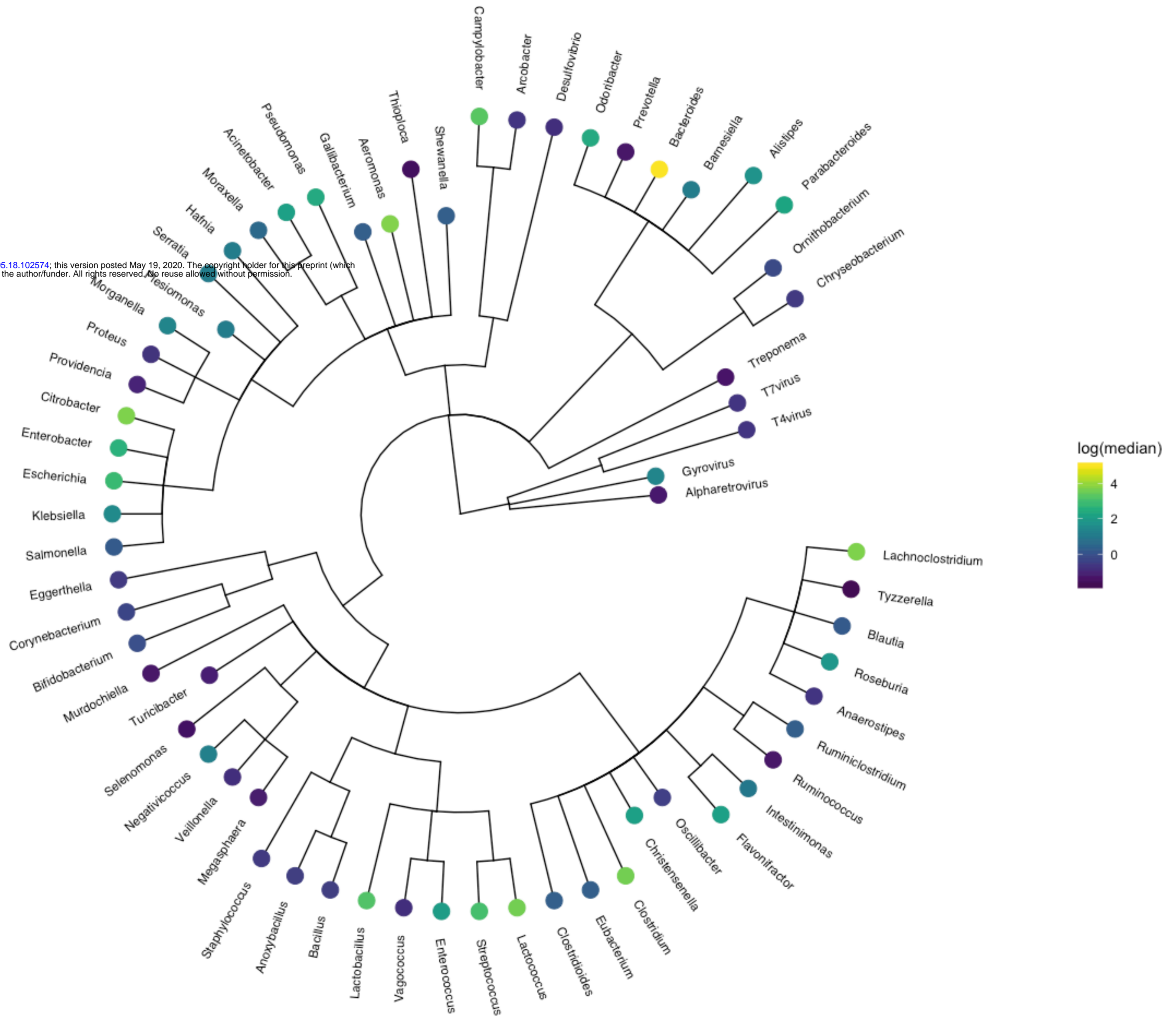


Beta diversity clustering by sample (Aitchison distance)



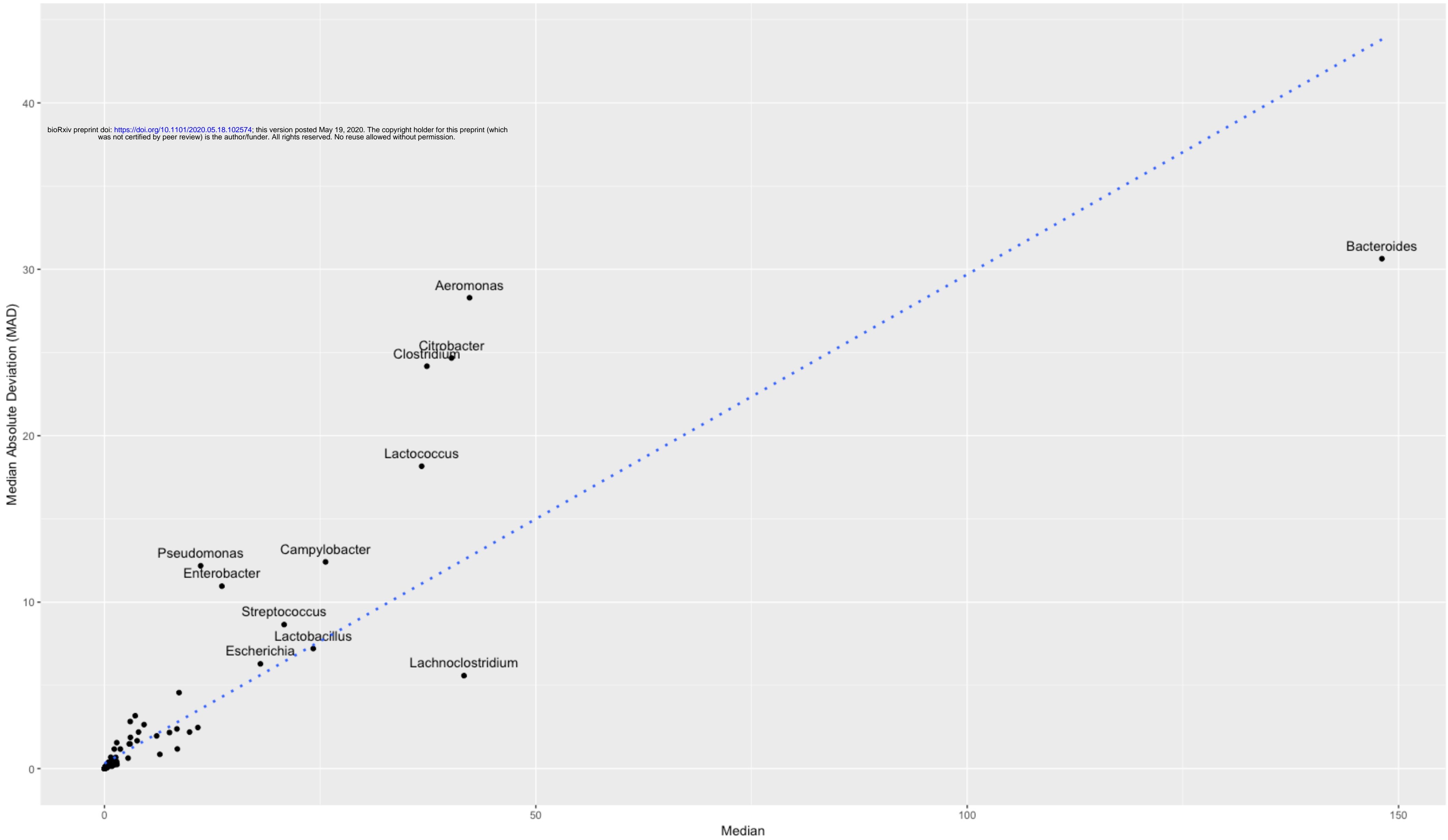
Phylogram of Consensus Microbial Genera by Abundance

bioRxiv preprint doi: <https://doi.org/10.1101/2020.05.18.102574>; this version posted May 19, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

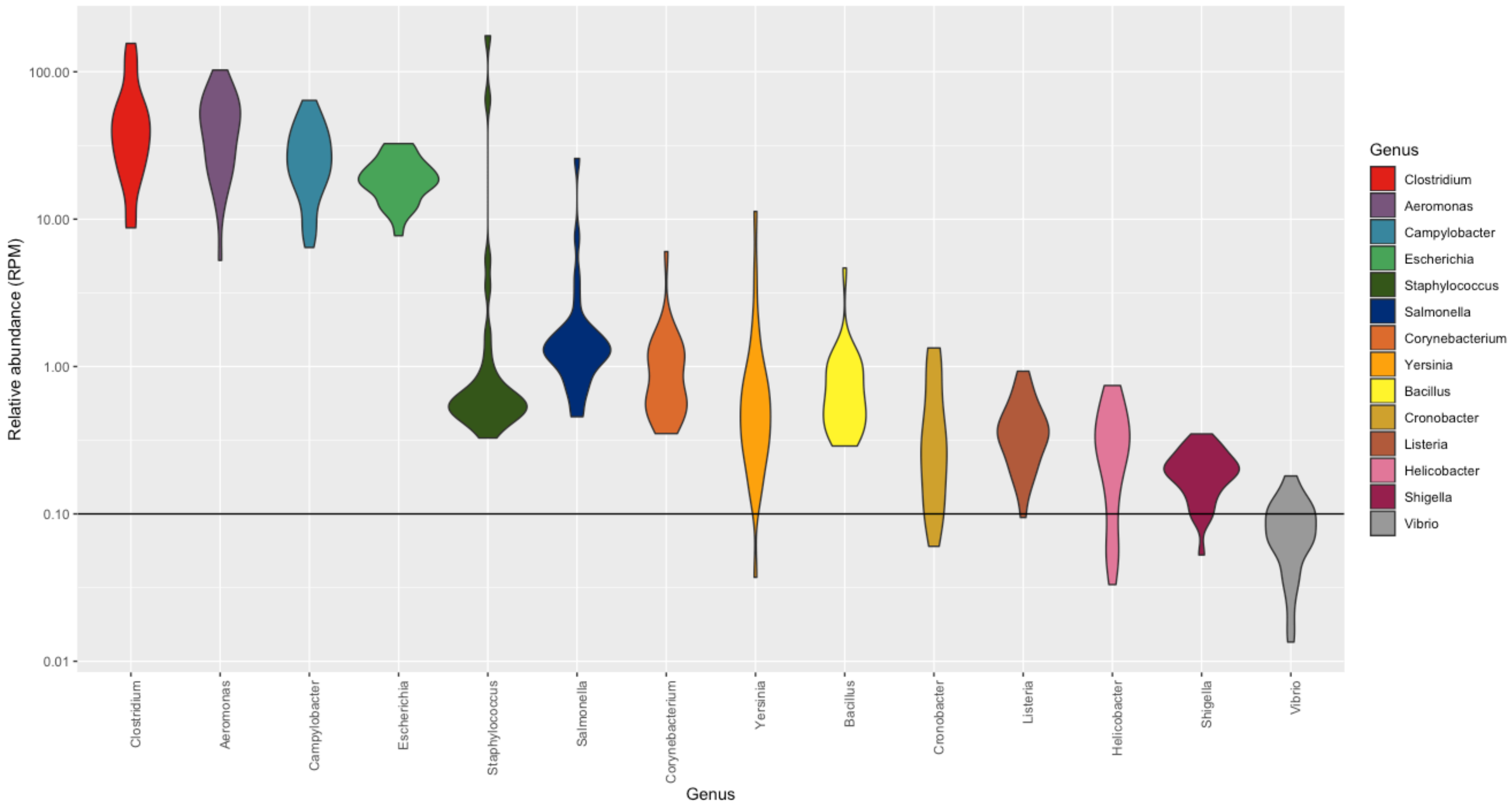


Median and MAD of microbial genera (RPM > 0.1)

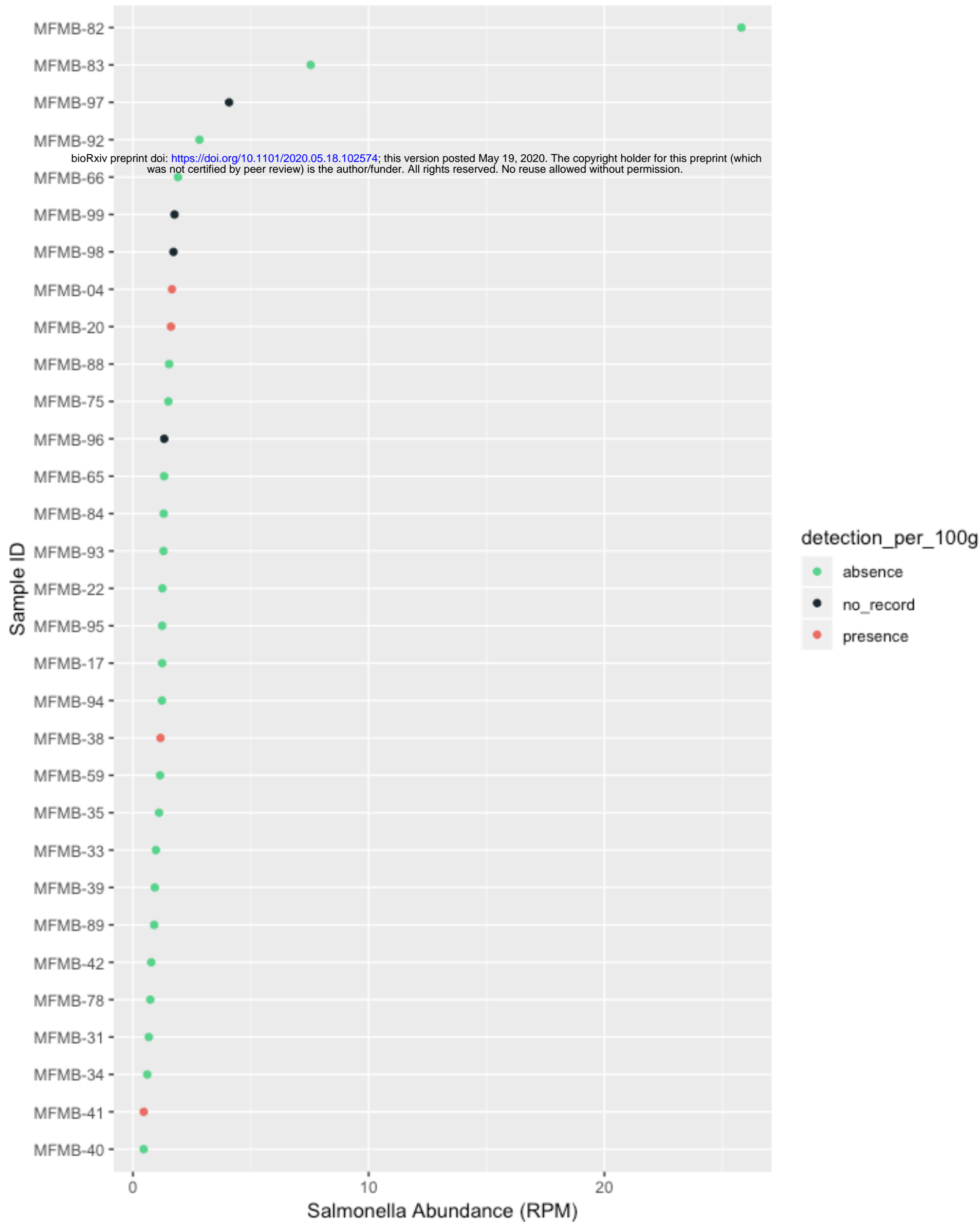
bioRxiv preprint doi: <https://doi.org/10.1101/2020.05.18.102574>; this version posted May 19, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.



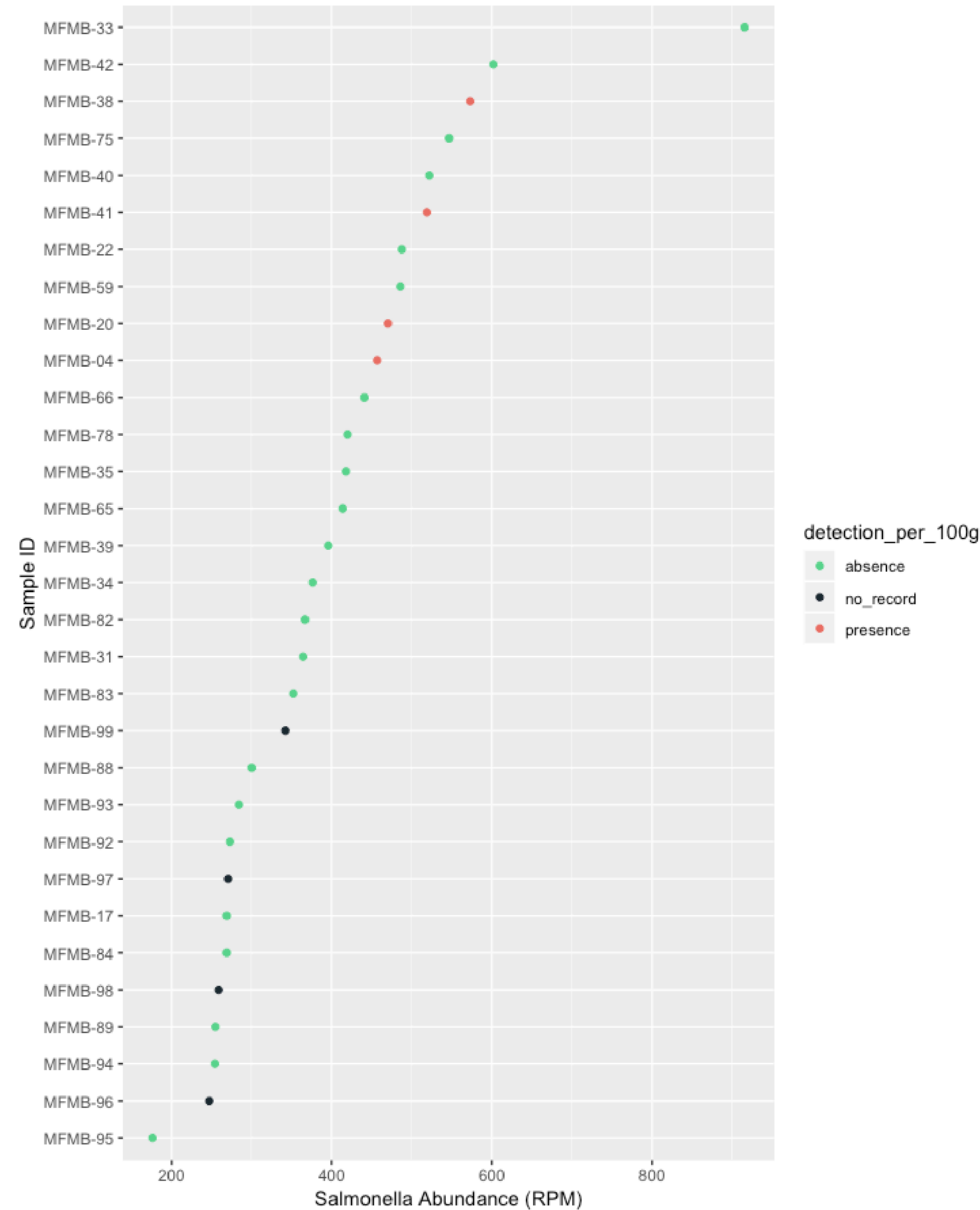
High priority food safety genera relative abundance



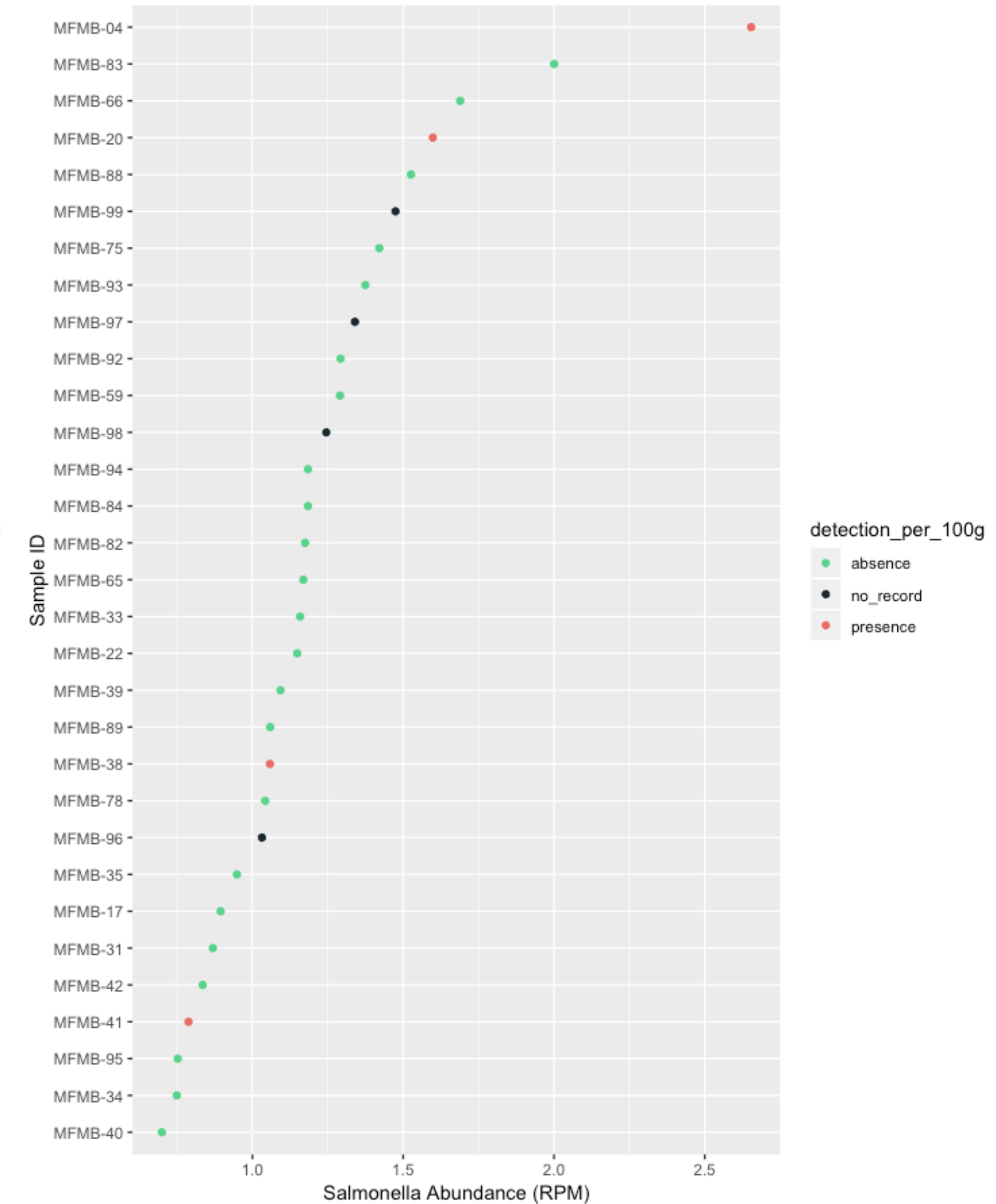
A. Salmonella Content from Metatranscriptome Classifications By Sample



B. Salmonella Alignments to Complete Genomes By Sample



C. Salmonella Alignments to ef-Tu By Sample



<i>Salmonella</i> -positive sample	k-mer Classification	Whole Genome Alignment	<i>ef-Tu</i> Alignment
MFMB-04	8th	10th	1st
MFMB-20	9th	9th	4th
MFMB-38	20th	3rd	21st
MFMB-41	30th	6th	28th
Rank sum test p-value	p=0.86	p=0.06	p=0.56

Bi-Serial Correlation (Supported Microbes & |Corr.| > 0.5)

