1 **FoodMicrobionet v4: a large, integrated, open and transparent database for food bacterial**
2 **communities.**
3
4 Eugenio Parente, Teresa Zotta, Annamaria Ricciardi
5
6 Scuola di Scienze Agrarie, Forestali, Alimentari ed Ambientali, Università degli Studi della
7 Basilicata, Potenza, Italy
8
9
10 **Corresponding author**
11
12 Prof. Eugenio Parente
13
14 Scuola di Scienze Agrarie, Forestali, Alimentari ed Ambientali,
15 Università degli Studi della Basilicata,
16 Viale dell'Ateneo Lucano 10, 85100 Potenza, Italy
17
18 Email: eugenio.parente@unibas.it
19 Tel.: +390971205561
20
21

22  **FoodMicrobionet v4: a large, integrated, open and transparent database for food bacterial**

23  **communities.**

24

25  Eugenio Parente, Teresa Zotta, Annamaria Ricciardi

26

27  **Abstract**

28  With the availability of high-throughput sequencing techniques our knowledge of the

29  structure and dynamics of food microbial communities has made a quantum leap. However,

30  this knowledge is dispersed in a large number of papers and hard data are only partly

31  available through powerful on-line databases and tools such as QIITA, MGnify and the

32  Integrated Microbial Next Generation Sequencing platform, whose annotation is not

33  optimized for foods.

34  Here, we present the 4$^{th}$ iteration of FoodMicrobionet, a database of the composition of

35  bacterial microbial communities of foods and food environments. With 180 studies and

36  10,151 samples belonging to 8 major food groups FoodMicrobionet 4.1.2 is arguably the

37  largest and best annotated database on food bacterial communities. This version includes

38  1,684 environmental samples and 8,467 food samples, belonging to 16 L1 categories and

39  196 L6 categories of the EFSA FoodEx2 classification and is approximately 4 times larger

40  than previous version (3.1, https://doi.org/10.1016/j.ijfoodmicro.2019.108249).

41  Using data in FoodMicrobionet we confirm that taxonomic assignment at the genus level

42  can be performed confidently for the majority of amplicon sequence variants using the most

43  commonly used 16S RNA gene target regions (V1-V3, V3-V4, V4), with best results with

44  higher quality sequences and longer fragment lengths, but that care should be exercised in

45  confirming the assignment at species level.

46    Both FoodMicrobionet and related data and software conform to FAIR (findable, accessible,

47    interoperable, reusable/reproducible) criteria for scientific data and software and are freely

48    available on public repositories (GitHub, Mendeley data).

49    Even if FoodMicrobionet does not have the sophistication of QIITA, IMNGS and MGnify, we

50    feel that this iteration, due to its size and diversity, provides a valuable asset for both the

51    scientific community and industrial and regulatory stakeholders.

52

53    **Key words**: amplicon targeted high-throughput sequencing; 16S metagenomics; bacterial
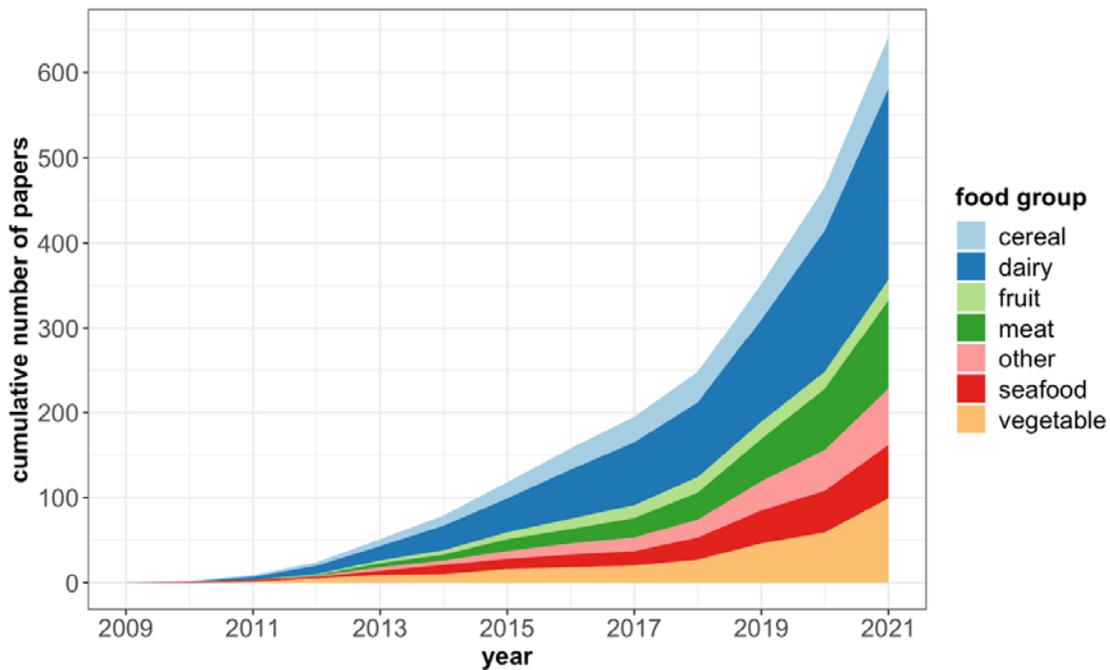
54    microbiota; database

55

56    **1. Introduction**.

57

58    In less than 15 years, next generation sequencing has added several layers of complexity to

59    the study of food microbiomes. At the time of writing this article (January 2022) searches

60    with the key words "microbiome" OR "microbiota" in Scopus and in Web of Science

61    retrieved 118,416 and 127,424 hits, respectively. The number of papers using amplicon

62    targeted next generation sequencing to study the structure of food bacterial communities

63    has been rising exponentially since the first pioneering publications (Humblot and Guyot,

64    2009; Roh et al., 2010) (see Figure 1).

65    The evolution of -omic approaches for the study of the food microbiome and the

66    methodological challenges in wet laboratory methods, sequencing and bioinformatic and

67    statistical approaches have been reviewed multiple times (Bokulich et al., 2020; De Filippis

68    et al., 2018; Pollock et al., 2018). While amplicon targeted approaches are still prevalent, the

69    number of studies using metagenomic or multi-omics approaches is rapidly increasing with

3

70    the decreasing cost of sequencing and the increase in power of sequencing platforms and

71    computational tools (Yap et al., 2021).



72

73    **Figure 1.** Cumulative number of papers using metataxonomic or metagenomic approaches

74    for the study of food microbial communities. The data derive from a manually curated

75    search on Web of Science[TM].

76

77    Shotgun metagenomics offers significantly higher resolution (down to the strain and

78    perhaps to single nucleotide polymorphism variants (Hildebrand, 2021) and paves the way

79    to accurate source tracking for contamination (De Filippis et al., 2020) and to a new

80    paradigm in food safety (Kovac, 2019). However, amplicon targeted approaches still have

81    the power to describe, in a sensitive and cost-effective way, the structure of food microbial

82    communities down to the genus level and possibly below (Callahan et al., 2016a; Johnson et

83    al., 2019).

84    More than 640 papers describing the food microbiota have been published since 2009.

85    Exploitation of this wealth of information for metastudies or for the development of

86    descriptive or predictive tools requires FAIR (findable, accessible, interoperable,

87    reusable/reproducible: Lamprecht et al., 2020; Wilkinson et al., 2016) data and software.

88    Three large on line repositories on microbiome data have appeared in recent years: IMNGS

89    (Lagkouvardos et al., 2016), QIITA (Gonzalez et al., 2018) and MGnify (Mitchell et al., 2019).

90    The Integrated Microbial Next Generation Sequencing platform directly accesses 16S rRNA

91    gene targeted next generation sequencing data in NCBI Short Read Archive and provides

92    powerful tools for searching for taxa or sequences or for performing analysis of user-

93    deposited 16S sequences (Lagkouvardos et al., 2016).

94    MGnify hosts metagenomic, metatranscriptomics and metataxonomic datasets which are

95    either contributed by users or obtained from public repositories and offers powerful

96    platforms for data analysis but does not allow the integration of data from different studies

97    (Mitchell et al., 2019).

98    QIITA offers a powerful suite of tools for the analysis of sequences of public and private

99    datasets, and allows the search and integration of data from different studies (Gonzalez et

100    al., 2018).

101    At the time of writing of this paper (January 2022) MGnify included 3,696 public studies, and

102    325,323 public samples, but only 83 studies/2,805 samples on food biomes

103    (https://www.ebi.ac.uk/metagenomics/) while QIITA included 620 public studies and

104    303,313 public samples but only a very limited number of public studies on food biomes

105    (https://qiita.ucsd.edu/stats/). The data and interfaces for all these tools respond well to

106    FAIR principles for scientific data (Wilkinson et al., 2016) but, unfortunately, the structure of

107    the metadata for studies and samples is not optimized for foods.

108    Some years ago, we created a database for metataxonomic studies on food bacterial

109    communities, FoodMicrobionet (Parente et al., 2016, 2019), whose main strength is the

110    annotation system for studies and samples, based on the hierarchical classification of foods

111    developed by the European Food Safety Authority, FoodEx 2 (E.F.S.A., 2015). Here we

112    describe the structure and new features of the latest version of the database. In addition,

113    we provide two proofs of concept on how the rich metadata structure of FoodMIcrobionet

114    can be used to demonstrate the effect of target region on the resolution of taxonomic

115    assignment of amplicon targeted metagenomics for food bacteria.
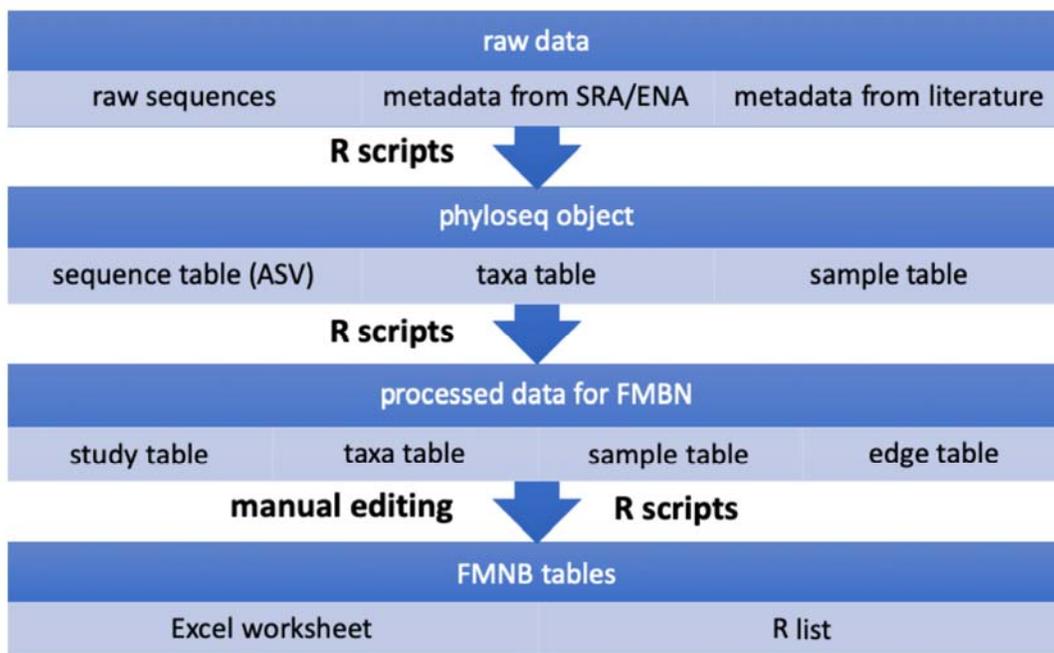
116

117    **2. Materials and methods.**

118

119    **2.1 Feeding data to FoodMicrobionet.**

120    The procedure used to add data to FoodMicrobionet has not changed since the last version

121    (Parente et al., 2019) and it is represented schematically in Figure 2.

122    With the exception of 33 studies, which were originally provided by research partners as

123    abundance tables for taxa and sample metadata tables (Parente et al., 2016), studies in

124    FoodMicrobionet are added to the database by reprocessing raw sequence data from NCBI

125    Short Read Archive (SRA), and by using metadata from SRA and from the scientific papers

126    for annotation.

127    Processing of sequences is carried out in R (R Core Team, 2021) using a modified version of

128    the Bioconductor pipeline for amplicon targeted sequence analysis, with DADA2 for

129    Amplicon Sequence Variant (ASV) inference and SILVA v138.1 for taxonomic assignment

130    (Callahan et al., 2016a; Callahan et al., 2016b). This results in the production of phyloseq

131    (McMurdie and Holmes, 2013) objects which are processed using R scripts and imported in

132    Microsoft Excel for further manual editing of study and sample metadata. Finally, a R script

133    is used to process Excel tables and for quality control checks. All scripts are publicly available

134    in the FoodMicrobionet GitHub repository (https://github.com/ep142/FoodMicrobionet).



135

136    **Figure 2.** Schematic workflow showing how raw sequences and their metadata and

137    additional information from the scientific literature are used to assemble data for
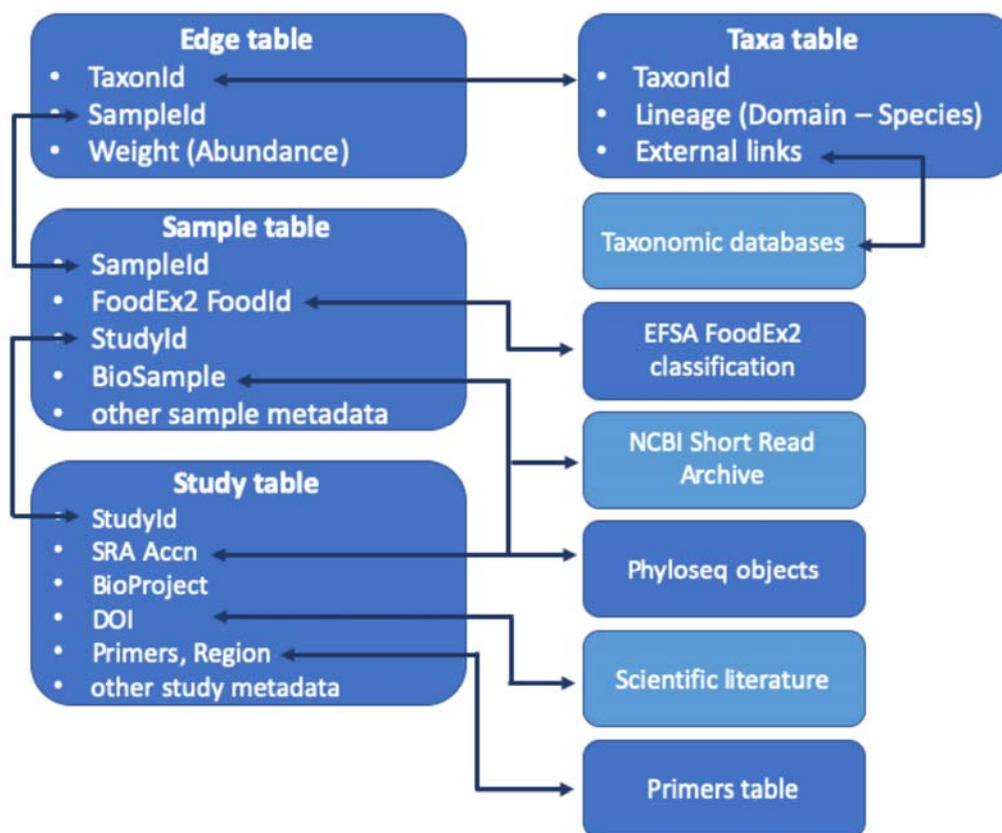
138    FoodMicrobionet.

139

140    **2.2 The structure of FoodMicrobionet v4.**

141    The structure of the database is schematically shown in Figure 3.

142    Studies, samples and taxa have all a rich metadata structure. The structure of the tables is

143    described in detail on GitHub

144    (https://github.com/ep142/FoodMicrobionet/blob/master/FoodMicrobionet_tablespecs.m

7

145    d), on Mendeley data (https://data.mendeley.com/datasets/8fwwjpm79y/5) and in

146    supplementary material.



147

148    **Figure 3.** Schematic representation of the structure of FoodMicrobionet tables (Study,

149    Primers, Samples, FoodEx2, Taxa, Edges) and their relationship with phyloseq (McMurdie

150    and Holmes, 2013) objects and external databases.

151

152    Additions in version 4 include four fields which complete information for bioinformatic

153    processing and one field for geolocation for studies, geolocation information for samples,

154    and two further reference tables (primers and the FoodEx2 Exposure Hierarchy classification

155    (E.F.S.A., 2015).

156    FoodMicrobionet is available either as R lists, which allow experienced programmers to run

157    their searches and analyses in the most flexible and sensitive way, and as an interactive

158    Shiny app (Parente et al., 2019). The latter requires minimum installation and configuration

159    and allows users to perform searches using a large number of criteria, to perform

160    aggregation of samples and taxa, to rapidly reach external resources using hyperlinks, to

161    export data in a variety of formats, and to obtain and save graphs and tables (Parente et al.,

162    2019).

163

164    **2.3 Proof of concept 1: on the taxonomic resolution of amplicon targeted metagenomics**

165    **for food bacterial communities**.

166    Using the metadata available in FoodMicrobionet we tabulated the frequency of taxonomic

167    level assignments at the species, genus, family, order, class and phylum level for studies 34

168    to 180 (i.e. all studies for which sequence processing had been done using the procedure

169    described in section 2.1). Graphs and tables were generated using a R script (ide_depth.R,

170    available on GitHub,

171    https://github.com/ep142/FoodMicrobionet/tree/master/the_real_thing/R_lists).

172

173    **2.4 Proof of concept 2: using ASV for in depth analysis of taxonomic assignments**

174    One of the major changes in FoodMicrobionet v4 is that ASV for each study can be directly

175    accessed using the phyloseq objects created with the pipeline described in section 2.1. This

176    in turn allows comparisons among ASV obtained in different studies using the same target

177    region, and with reference sequences. To demonstrate this, we wrote a script which

178    performed the following actions:

179 1. Search the database for two genera including pathogenic species (*Listeria, and*

180   *Salmonella*) and identify the samples and studies in which they occur

181 2. Create graphs and tables showing their prevalence and abundance

182 3. Use study and sample accession numbers to retrieve ASV sequences from the

183   phyloseq objects

184 4. Divide the sequences in groups (depending on the 16S RNA gene target region), carry

185   out taxonomic assignment using RDP v18,  and compare the sequences with

186   reference sequences and outgroups extracted from the SILVA v138.1 reference

187   database (two randomly extracted reference sequences for each species were used)

188 5. For each group, perform alignment and estimate Maximum Likelihood phylogenetic

189   trees using the procedure described in Callahan et al. (2016b)

190 6. Annotate and plot phylogenetic trees using *treeio, tidytree* (Wang et al., 2020),  and

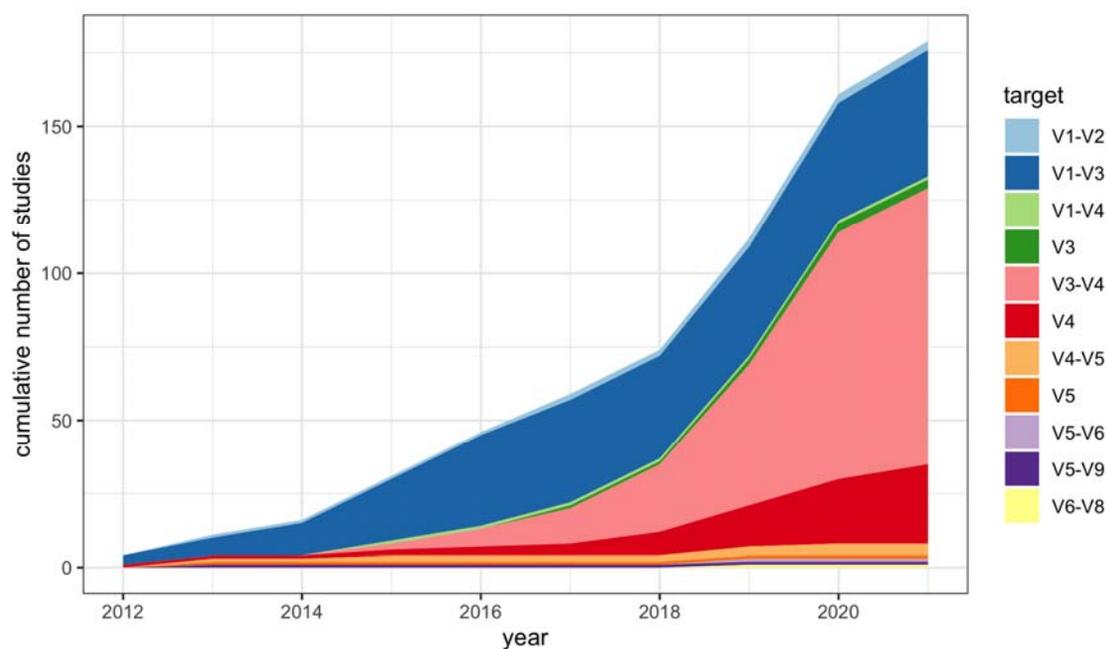191   *ggtree* (Yu, 2020) R packages

192

193 **3. Results and discussion.**

194

195 **3.1 FoodMicrobionet facts and figures.**

196 FoodMicrobionet has grown significantly since its last iteration (Parente et al., 2019; Figure

197 4): the number of studies in version 4.1.2 has grown from 44 to 180, and the number of

198 samples from 2,234 to 10,151.

199 The growth of FoodMicrobionet reflects the growth of published studies and the

200 distribution of target gene regions reflects the shift in technology in Next Generation

201 Sequencing, from the defunct Roche 454 platform (with most studies targeting the V1-V3

10

202    region) to Illumina platforms, with V4 and V3-V4 regions of the 16S rRNA gene as the most

203    frequent targets (Figure 4, Supplementary Table 1).



204

205    **Figure 4.** Cumulative distribution of studies in FoodMicrobionet 4.1.2, by year and by 16S

206    rRNA gene target region.

207

208    The number of reads for sample after processing is also saved in the database

209    (Supplementary Figure 1): 70% samples have >$10^4$ reads after processing, but a significant

210    number of samples in studies targeting the V4 region has >$5 \times 10^4$ reads. This makes the

211    detection of rare components of the microbiota possible.

212    The variety of food groups and food types in FoodMicrobionet is also very large: samples in

213    FMBN belong to 16 major food groups (L1 level of FoodEx2 exposure classification;

214    Supplementary  Table 2). Approximately 80% of samples belong to the first three categories

215    (Milk and dairy products, Meat and meat products, Vegetables and vegetable products).

216    Samples in FMBN are further classified using levels L4 and L6 of the FoodEx2 exposure

217    classification, and additional fields (which allow to identify raw products, intermediates or

218    finished products, the level of thermal treatment and the occurrence of spoilage and/or

219    fermentation) allow a finer classification. Samples in FMBN belong to 109 L4 food groups

220    and 197 L6 food groups. There are 163 individual food types, and, if further information on

221    samples (nature, heat treatment, spoilage/fermentation) are used, there are 316 unique

222    combinations (Parente et al., 2019). This number and variety of samples makes

223    FoodMicrobionet the largest database of metataxonomic data on food bacterial

224    communities, with significantly more samples compared to QIITA and MGnify.

225    FoodMicrobionet stores samples from 51 countries, but the 90% of samples are from 14

226    countries (Supplementary Figure 2). This does not reflect the distribution of samples and

227    studies in published studies but, rather, the distribution of those which are available in NCBI

228    SRA.

229    There are currently 9,098 taxa in the taxa table of FoodMicrobionet. Taxa have a unique

230    numeric and text identifier and may represent ASV identified at the species (4,497, 49.4% of

231    taxa) or genus level (3,259, 35.8%) or above using the DADA2 assignTaxnomy() and

232    assignSpecies() functions. The % of the taxa and sequences identified at the genus level or

233    below varies by study, depending on the quality and length of sequences and on the gene

234    target. Length of reads (in bp) in FoodMicrobionet studies varies between 150 and 610 bp

235    (median 422).

236    FoodMicrobionet is fully connected to external databases: external links (as dynamically

237    built Uniform Resource Locators) in the study, sample and taxa tables allow rapid access to

238    external taxonomic databases (NCBI taxonomy, the List of Prokaryotic Names with Standing

239    in Nomenclature and the Florilege database, Falentin et al., 2017

240    http://migale.jouy.inra.fr/Florilege/#&about) and to the scientific literature (via DOI). Using

12

241     NCBI SRA Study accession number it is possible to access fine grained data on ASV sequence

242     and abundance stored in the phyloseq objects obtained from processing the raw sequences

243     (these are not public but are available on request).

244

245     **3.2 Is FoodMicrobionet FAIR?**

246     FoodMicrobionet data and software are free (both are under MIT licence

247     https://opensource.org/licenses/MIT), open and highly reusable and support analysis

248     protocols which are reproducible. We have done our best to conform as closely as possible

249     to criteria for FAIR (findable, accessible, interoperable, reusable/reproducible) data and

250     software sharing (Lamprecht et al., 2020; Wilkinson et al., 2016).

251     Both the database and the software are findable (using searches on Google, Mendeley data,

252     or GitHub, for example) and deposited on permanent repositories (Mendeley data, GitHub)

253     with unique identifiers. Since the database is not available on line except in the form of R

254     lists or Excel files, data within FoodMicrobionet may not be directly findable by automated

255     machine searches (and insofar they are not machine operable); however, the wealth of

256     contextual metadata for all the main tables (studies, samples, taxa) makes it possible to

257     devise precisely targeted searches.

258     FoodMicrobionet is accessible through the above-mentioned repositories and through our

259     website, and we are confident that enough metadata are provided in these repositories to

260     easily reach the resources. In terms of user accessibility (which is not a criterion in FAIR

261     principles), we have done our best to make the resource accessible to both expert (the

262     database in the form of a R list can be used for fine-tuned searches and analyses using R

263     scripts) and moderately expert users. The latter can, with a minimum effort, download and

264     install the R Shiny app, ShinyFMBN (Parente et al., 2019), which, once launched in R, allows

265    easy access via an interactive and intuitive interface. A detailed manual for the app is

266    available on Mendeley data (https://data.mendeley.com/datasets/8fwwjpm79y/4).

267    Although the app could be easily deployed on RStudio Shiny apps server

268    (https://www.shinyapps.io), we feel that this would make the use unnecessarily slow.

269    FoodMicrobionet is fully interoperable. The sample classification is based on a robust

270    hierarchical classification, FoodEx2, rather than on arbitrary keywords, and dynamic links

271    are created in the studies, samples and taxa tables to reach a number of other databases.

272    Conversely, other databases like Florilege might quite simply create new accessions using

273    FoodMicrobionet and metadata in FMBN can, in principle, be used to populate QIITA and

274    MGnify, by linking studies and samples via the SRA accession numbers.

275    Data and products of search results are highly reusable, for the same reasons. In addition,

276    the objects exported by the app are in formats which are compatible with metataxonomic

277    analysis pipelines (like phyloseq and ShinyPhyloseq: (McMurdie and Holmes, 2013, 2015);

278    MicrobiomeAnalyst: (Chong et al., 2020; Dhariwal et al., 2017); CoNet: (Faust and Raes,

279    2016); graph visualization and analysis software like Cytoscape and iGraph: (Csardi et al.,

280    2006; Shannon et al., 2003), microbial association network inference tools: (Kurtz et al.,

281    2015; Peschel et al., 2021).

282    Use cases and example workflows have been illustrated in a previous work (Parente et al.,

283    2019) and we have tried to demonstrate this approach in a series of proof-of-concept

284    metastudies (Parente et al., 2020, 2021; Zotta et al., 2021). The code for generating graphs

285    and statistical analyses is fully reproducible and reusable

286    (https://data.mendeley.com/datasets/8fwwjpm79y/4;

287    https://github.com/ep142/MAN_in_cheese) and allows to reproduce the figures and tables

288    whenever a new version of FoodMicrobionet is published.
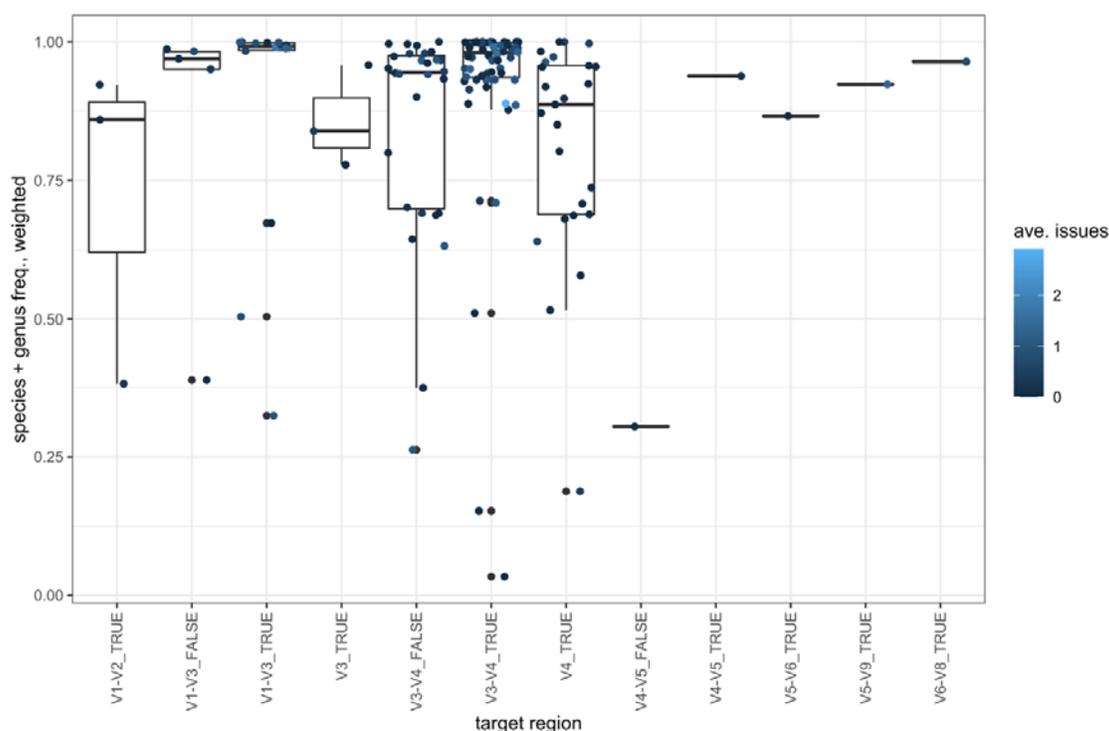
14

289

**3.3 Proof of concept 1: on the taxonomic resolution of amplicon targeted metagenomics**

**for food bacterial communities.**

The procedure for adding data to FMBN involves, starting from version 3 (studies 34 to 180)

the use of a modification of the BioConductor workflow for microbiome data analysis

(Callahan et al., 2016b). This procedure, which infers Amplicon Sequence Variants using

DADA2, has performed well in benchmarking and can be used to compare data across

multiple studies (Callahan et al., 2016a. Callahan et al., 2017). ASV inference is increasingly

been used in the study of food bacterial communities (55% of studies 160 to 180 in

FoodMicrobionet originally used DADA2; see Supplementary Table 3 for references). The

ability to perform taxonomic assignment to the lowest possible level (species) is clearly

important in food microbial ecology, because different members a given genus may have

different roles in foods (beneficial, detrimental, pathogenic); this is for example the case for

*Clostridium, Bacillus, Staphylococcus, Corynebacterium,* and many other genera of food

related bacteria. Different variable regions of the 16S RNA gene have been historically used

as targets for amplicon targeted metagenomics, and FoodMicrobionet provides a

comprehensive sampling of the targets used (Figure 4).

The median and 90° percentile values for the frequency of taxonomic assignment at the

genus level or below in FMBN studies from 34 to 180 is shown in Supplementary Table 4.

The weighted (by sequence abundance, Figure 5) and unweighted (Supplementary Figure 3)

frequency of assignment at the genus level or below varied between studies, even within a

given target region (Figure 5). This was apparently related to sequence length and target

region, and a clear relationship was found with at least one indicator of average within-

study sequence quality. In fact, for paired end sequences for the longest target regions,

15

313     merging of paired ends was not always possible due to bad quality of sequences toward the

314     5' end (data not shown) and this prevented the overlap of forward and reverse sequences.

315     For these sequences, the BioConductor workflow used in FoodMicrobionet allows

316     taxonomic assignment down to the genus level, but not down to the species level (Callahan

317     et al., 2016b).



318

319     **Figure 5.** Box and jitter plots showing the weighted (by sequence relative abundance)

320     distribution of frequencies of taxonomic assignments at the genus level or below in

321     FoodMicrobionet studies 34 to 180. The average values for the number of issues

322     encountered during bioinformatic processing (high sequence losses during filtering or

323     chimera removal, low number of final sequences, low diversity) is also shown.

324

325     However, no clear relationship was found with the other indicator of sequence quality

326     provided by FoodMicrobionet, i.e. the average number of issues during bioinformatic

327    processing, see table specifications in Supplementary material). Overall, the median value of

328    the frequency of taxonomic assignment at the genus level or below ranged from 0.640 and

329    0.898, with the lowest values for the shortest regions and studies with the worse sequence

330    quality. However, when the number of sequences for each ASV is taken into account these

331    figures may change significantly, and median values for taxonomic assignment at the genus

332    level or below as high as 0.98 (overlapping V3-V4 region) or 0.99 (overlapping V1-V3 region)

333    can be obtained (Figure 5). Shorter regions still provide a reasonably good performance

334    (with weighted median frequencies of genus assignments of 0.73 and 0.80 for V3 and V4

335    respectively). However, species assignments were much less frequent, with median values

336    for weighted frequencies ranging from 0.025 to 0.381 (Supplementary Table 5). Median

337    weighted values for species level assignment were 0.38 and 0.30 for regions V1-V3 and V3-

338    V4, respectively, but as low as 0.21 for V4, a frequently used target in large recent studies.

339    Differences in the frequencies of taxonomic assignment were also observed for different

340    phyla. This is illustrated for the four most abundant phyla in FoodMicrobionet (*Firmicutes,*

341    *Proteobacteria, Actinobacterota, Bacteroidota*; SILVA taxonomy is used for higher taxa) in

342    Supplementary Figures 4 and 5. For some targets the ability to perform taxonomic

343    assignment at the genus level or below was clearly lower and/or more variable.

344    These results are in good agreement with a recent study which compared taxonomic

345    assignment for different target regions within the 16S RNA gene (Johnson et al., 2019). The

346    possibility of assigning taxonomy down to the species level was found to differ among

347    regions, with species level assignments are significantly less frequent for shorter regions

348    compared to the full 16S RNA gene, which can now be sequenced using 3[rd] generation High

349    Throughput Sequencing platforms, like PacBio and Oxford Nanopore (Johnson et al., 2019).

350    In addition, differences in the ability to perform taxonomic assignment at the genus and

351    species level for different phyla may also explain why the weighted and unweighted

352    frequencies of identification differ: abundant sequences often belong to taxa which are well

353    represented in taxonomic reference databases (data not shown).

354    In FoodMicrobionet, target region and composition of the microbiota of the study are

355    clearly not independent. Given the metadata structure in FoodMicrobionet one could, at

356    least in principle, compare the taxonomic assignment for different food groups (which

357    might, in turn, reflect differences in microbial community composition), but this might result

358    in too many different combinations. However, at least for the target regions for which a

359    large number of studies is available (V3-V4, V4, and, to a lesser extent, V1-V3) we feel that

360    the results offer a wide enough coverage of food groups and clearly indicate that for

361    *Actinobacterota* and *Bacteroidota* the pipeline used in FoodMicrobionet may offer a lower

362    degree of success in taxonomic assignment compared to *Proteobacteria* and *Firmicutes*.
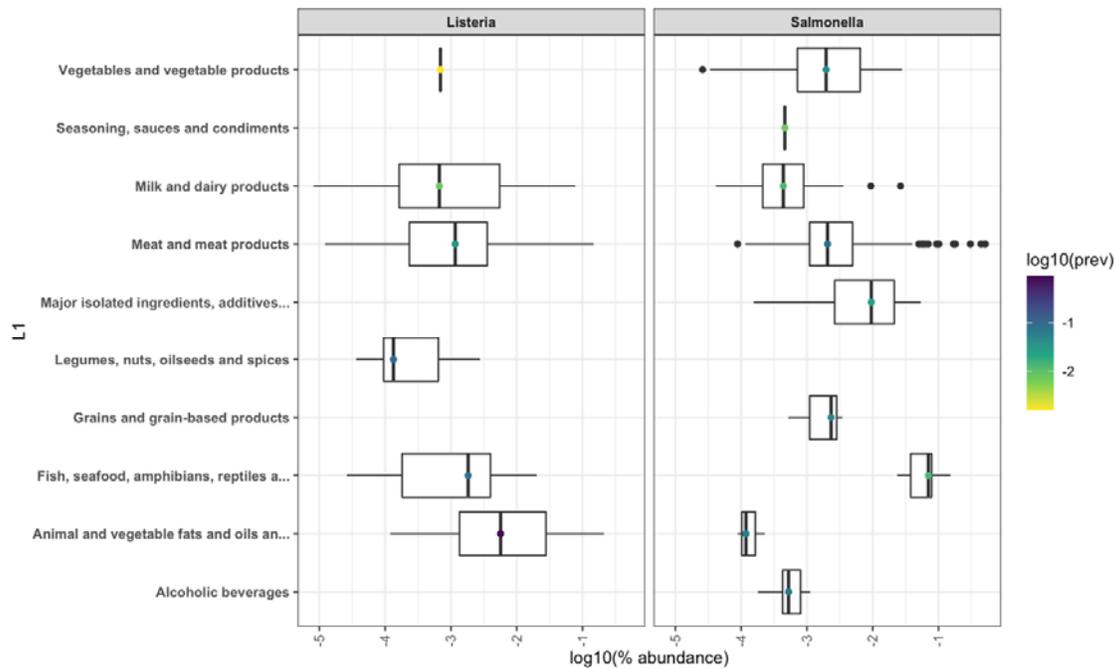
363

364    **3.4 Proof of concept 2: using Amplicon Sequence Variants for in depth analysis of**

365    **taxonomic assignments**

366    To further demonstrate the need to exert caution in taxonomic assignment down to the

367    species level with relatively short targets we developed a second proof of concept. First,

368    we searched the database for samples containing members of genera *Listeria* and

369    *Salmonella*. While metagenomic approaches, due to their high resolution, are certainly

370    more useful than metataxonomic approaches in food safety studies (Cocolin et al., 2018;

371    Jagadesaan et al., 2019; Kovac, 2019), the latter may still have value for studying the

372    microbial ecology of food borne pathogens. However, this requires taxonomic assignment

373    to a level which is low enough to discriminate species, or species groups, which are relevant

374    for human or animal health. The distribution of the abundance of members of the genera

18

375     *Listeria* and *Salmonella* in L1 food categories of the EFSA FoodxEx2 classification is shown in

376     Figure 6.



377

378     **Figure 6.** Distribution of the abundance of six genera including pathogenic bacteria in

379     FoodMicrobionet samples. The L1 level of food classification of the EFSA FoodEx2

380     classification is shown. Prevalence is shown as a colour scale. Environmental samples were

381     excluded from the analysis.

382

383     With some exceptions in which abundance was >0.1% and as high as 0.5% of total

384     sequences (Supplementary Table 6), these genera occur with a low prevalence and

385     abundance (typically <0.01%) and would normally be discarded by abundance and

386     prevalence filters which are normally applied when processing microbiome data (Callahan et
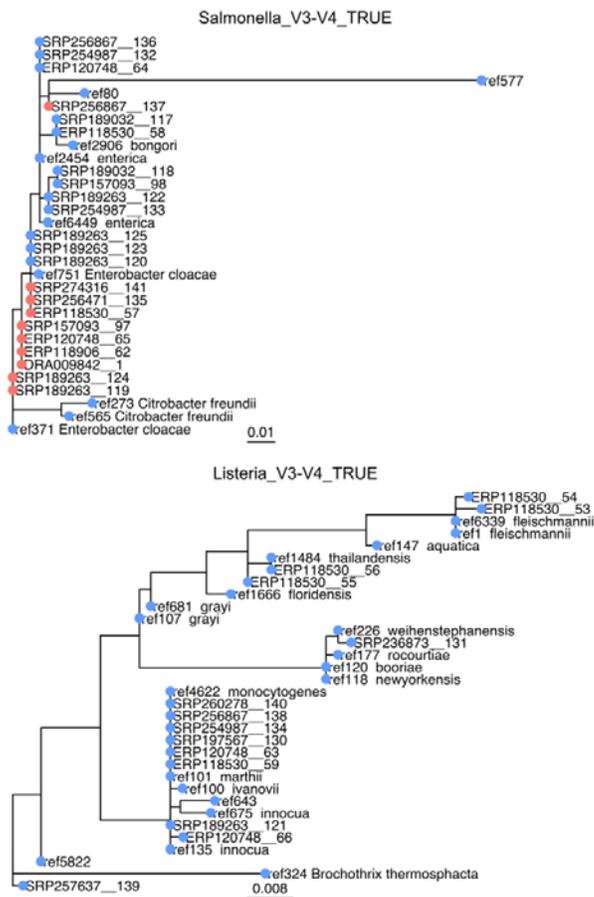
387     al., 2016b). The low abundance and prevalence and the occurrence of some genera in

388     unexpected environments (like *Salmonella* in alcoholic beverages) may rise the doubt that

389     their detection is due to contamination or to errors in sequence processing or taxonomic

390    assignment. Unfortunately, although it is well known that contamination may severely

391    affect the results of microbiome analysis, especially in low biomass samples (Dahlberg et al.,

392    2019; Davis et al., 2018; Pollock et al., 2018), the use of blanks and control and the

393    application of statistical procedures for the removal of contamination (Davis et al., 2018) is

394    very rare. Moreover, even if mock communities may assist in benchmarking the

395    bioinformatic processing of sequences in metataxonomic studies (Bokulich et al., 2020;

396    Pollock et al., 2018) they, too, are very rarely used in food microbial ecology studies.

397    While providing conclusive results on the occurrence of contamination or on the quality of

398    bioinformatic processing is impossible, the accessibility of ASV in FoodMicrobionet using

399    study and sample accession numbers provides an opportunity for checking the quality of the

400    taxonomic assignment and carry out direct comparison with reference sequences.

401    We therefore re-identified all sequences belonging to *Listeria* and *Salmonella* using the RDP

402    v18 trainset reference database. The results are shown in Supplementary Figure 6. While for

403    *Listeria* the two databases resulted in matching identifications, for *Salmonella*, a high

404    proportion of sequences were assigned to other genera (*Enterobacter, Citrobacter*) using

405    RDP. Differences in taxonomic assignment due to the use of different reference databases

406    are not surprising (Ramakodi, 2022; Werner et al., 2012) even if SILVA and RDP often

407    produce matching assignments (Smith et al., 2020).

408    Finally, we generated phylogenetic trees for sequences grouped by region, and included

409    reference sequences from the SILVA v138.1 taxonomic reference database, including

410    outgroups (*Brochothrix thermosphacta* for *Listeria* and *Citrobacter freundii* for *Salmonella*).

411    Reference sequences for *Enterobacter cloacae* were also included for *Salmonella*. Results for

412    overlapping paired sequences for region V3-V4 are shown in Figure 7, while combined

413    results for the V3-V4 and V4 region are shown in Supplementary Figure 7. Results with other

20

414     regions were similar (data not shown). Classification obtained by probabilistic assignment

415     (as in the naïve Bayesian classifier, Wang, 2007) and phylogenetic tree inferences are based

416     on different approaches and are not easy to compare.



417

418     **Figure 7.** Maximum likelihood phylogenetic trees for amplicon sequence variants (ASVs) for

419     overlapping paired end sequences for region V3-V4 of the 16S RNA gene identified as

420     *Salmonella* or *Listeria*. ASVs are identified by the accession number of the study to which

421     they belong and by a random progressive integer. Reference sequences extracted from

422     SILVA v138.1 taxonomic reference database are also included. Colored dots indicate

423     sequences for which taxonomic assignment with SILVA v138.1 and RDP trainset 18 matched

424     (blue) or not (red).

425

426     However, for *Listeria,* reference sequences for different species grouped in several clades

427     (with slight differences in grouping for different regions). Due to small phylogenetic distance

428     within each clade, with reference sequences differentiated by a very small number of

429     nucleotide changes, it is not surprising that species assignment by the naïve Bayesian

430     classifier was often not successful. A single ASV (SRP257637_139) did not group with *Listeria*

431     reference sequences. As to *Salmonella*, the majority of ASVs grouped with *Salmonella*

432     reference sequences, even when taxonomic assignment at the species level with RDP was

433     different. However, at least for the V3-V4 region, one reference sequence for *Enterobacter*

434     *cloacae* clustered with *Salmonella*. This may explain differences in assignment at the genus

435     level with the two databases.

436     It is well known that accuracy of taxonomic assignment with the naïve Bayesian classifier

437     may vary for different regions (Wang, 2007) and that sequences in reference databases may

438     have erroneous taxonomic annotations (Pollock et al., 2018). However, we feel that our

439     results confirms that the ability to perform taxonomic assignments varies with different

440     regions of the 16S RNA gene (Johnson et al., 2019), and when using short fragments, even

441     when a taxonomic assignment at the species level is obtained, one should be wary of the

442     results.

443     Although one may question the value of taxonomic assignment at the species level,

444     especially for the shortest reads (Callahan et al., 2016a; Edgar, 2018; Johnson et al., 2019;

445     Meola et al., 2019; Pollock et al., 2018), due to the detailed information provided for both

446     studies and samples (gene target and region, the number of issues observed during

447     processing of raw sequences) and to the possibility of accessing to the processed sequences

448     (ASV) in phyloseq objects, users of FoodMicrobionet can make informed decision on how

449     and when taxonomic units should be combined at a level higher than the species, an

450     operation which is easily performed with the ShinyFMBN app (Parente et al., 2019).

451

452     **4. Conclusions.**

453     Even if FoodMicrobionet does not have the sophistication of QIITA, IMNGS and Mgnify, we

454     feel that this iteration, due to its size and diversity, provides a significant resource for both

455     the scientific community and industrial and regulatory stakeholders. Scientists can access

456     and use a variety of stand-alone or online software tools and ShinyFMBN to compare their

457     own results with literature results, carry out metastudies to answer a variety of scientific

458     questions, build reproducible analysis workflows, get quantitative data on the

459     ecology/distribution of bacteria of interest, use the database as an entry point for further

460     searches in other databases. The size of this version, which includes >9x10$^5$ taxon/sample

461     relationships, might even allow the machine learning approaches to predict contamination

462     patterns of food. The ability to rapidly retrieve information on prevalence/abundance of

463     taxon in different foods and on the structure of microbial communities in different food

464     types may be useful to both the industry and regulatory agencies. Information on the

465     distribution of beneficial genera and, to a lesser extent, species may find use for regulatory

466     purposes (for example to facilitate studies on the distribution of beneficial microorganisms

467     to evaluate their inclusion in the Qualified Presumption of Safety). The fine-grained data on

468     the structure of microbial communities for a large variety of raw materials, foods, food

469     environments may be useful for both process and product development purposes to

470     identify spoilage or contamination patterns, or for the design of microbiome-based starters.

471    We are committed to keep adding data to FoodMicrobionet, but the openness and

472    transparency of its software and documentation allows any interested party to create new

473    versions of the database or to significantly improve its structure and functionality.

474

475    **CrediT author statement**

476    **Eugenio Parente:** Conceptualization, Methodology, Software, Writing- Original draft

477    preparation. **Annamaria Ricciardi** Data curation, Writing – Reviewing and Editing. **Teresa**

478    **Zotta**: Data curation, Writing – Reviewing and Editing.

479

480    **Data statement.**

481    The database and related scripts and apps are available on GitHub

482    (https://github.com/ep142/FoodMicrobionet) and on Mendeley data

483    (https://data.mendeley.com/datasets/8fwwjpm79y/6). Phyloseq objects are available upon

484    request.

485

486    **Acknowledgements.**

487    This research did not receive any specific grant from funding agencies in the public,

488    commercial, or not-for-profit sectors.

489

490    **References**

491    Bokulich, N.A., Ziemski, M., Robeson, M.S., Kaehler, B.D., 2020. Measuring the microbiome:

492        best practices for developing and benchmarking microbiomics methods. Comput. Struct.

493        Biotechnol. J. 18, 4048–4062. https://doi.org/10.1016/j.csbj.2020.11.049

24

494    Callahan, B. J., McMurdie, P.J., Rosen, M.J., Han, A.W., Johnson, A.J.A., Holmes, S.P., 2016a.

495    DADA2: High-resolution sample inference from Illumina amplicon data. Nat. Met. 13,

496    581–583. https://doi.org/10.1038/nmeth.3869

497    Callahan, B.J., McMurdie, P.J., Holmes, S.P., 2017. Exact sequence variants should replace

498    operational taxonomic units in marker-gene data analysis. ISME J. 11, 2639–2643.

499    https://doi.org/10.1038/ismej.2017.119

500    Callahan, B. J., Sankaran, K., Fukuyama, J.A., McMurdie, P.J., Holmes, S.P., 2016b.

501    Bioconductor workflow for microbiome data analysis: from raw reads to community

502    analyses. F1000 Res. 5, 1492. https://doi.org/10.12688/f1000research.8986.2

503    Chong, J., Liu, P., Zhou, G., Xia, J., 2020. Using MicrobiomeAnalyst for comprehensive

504    statistical, functional, and meta-analysis of microbiome data. Nat. Protoc. 1–23.

505    https://doi.org/10.1038/s41596-019-0264-1

506    Cocolin, L., Mataragas, M., Bourdichon, F., Doulgeraki, A., Pilet, M.-F., Jagadeesan, B.,

507    Rantsiou, K., Phister, T., 2018. Next generation microbiological risk assessment meta-

508    omics: The next need for integration. Int. J. Food Microbiol. 287, 10–17.

509    https://doi.org/10.1016/j.ijfoodmicro.2017.11.008

510    Csardi, G., Nepusz, T., 2006. The igraph software package for complex network research,

511    InterJournal, Complex Systems 1695. https://igraph.org

512    Dahlberg, J., Sun, L., Waller, K.P., Östensson, K., McGuire, M., Agenäs, S., Dicksved, J., 2019.

513    Microbiota data from low biomass milk samples is markedly affected by laboratory and

514    reagent contamination. PLoS One 14, e0218257.

515    https://doi.org/10.1371/journal.pone.0218257

516    Davis, N.M., Proctor, D.M., Holmes, S.P., Relman, D.A., Callahan, B.J., 2018. Simple statistical

517        identification and removal of contaminant sequences in marker-gene and metagenomics

518        data. Microbiome 6, 226. https://doi.org/10.1186/s40168-018-0605-2

519    De Filippis, F., Parente, E., Ercolini, D., 2018. Recent past, present, and future of the food

520        microbiome. Annu. Rev. Food Sci. Technol. 9, 589–608. https://doi.org/10.1146/annurev-

521        food-030117-012312

522    De Filippis, F., Valentino, V., Alvarez-Ordóñez, A., Cotter, P.D., Ercolini, D., 2020.

523        Environmental microbiome mapping as a strategy to improve quality and safety in the

524        food industry. Curr. Opin. Food Sci. 38, 168–176.

525        https://doi.org/10.1016/j.cofs.2020.11.012

526    Dhariwal, A., Chong, J., Habib, S., King, I.L., Agellon, L.B., Xia, J., 2017. MicrobiomeAnalyst: a

527        web-based tool for comprehensive statistical, visual and meta-analysis of microbiome

528        data. Nucl. Ac. Res. 45, W180–W188. https://doi.org/10.1093/nar/gkx295

529    Edgar, R.C., 2018. Accuracy of taxonomy prediction for 16S rRNA and fungal ITS sequences.

530        PeerJ 6, e4652. https://doi.org/10.7717/peerj.4652

531    E.F.S.A., 2015. The food classification and description system FoodEx 2 (revision 2). EFSA

532        Supporting Publications 1–90. https://doi.org/10.2903/sp.efsa.2015.en-804

533    Falentin, H., Chaix, E., Dérozier, S., Weber, M., Buchin, S., Dridi, B., Deutsch, S.-M., Valence-

534        Bertel, F., Casaregola, S., Renault, P., Champomier-Verges, M.-C., Thierry, A., Zagorec, M.,

535        Irlinger, F., Delbes, C., Aubin, S., Bessières, P., Loux, V., Bossy, R., Dibie, J., Sicard, D.,

536        Nédellec, C. (2017, October). Florilege: a database gathering microbial phenotypes of

537        food interest. In Proceedings of the 4th International Conference on Microbial Diversity

538        2017, Bari, ITA (2017-10-24 – 2017-10-26). http://migale.jouy.inra.fr/Florilege/#&about

539   Faust, K., Raes, J., 2016. CoNet app: inference of biological association networks using

540      Cytoscape. F1000 Res. 5, 1519–14. https://doi.org/10.12688/f1000research.9050.2

541   Gonzalez, A., Navas-Molina, J.A., Kosciolek, T., McDonald, D., Vázquez-Baeza, Y., Ackermann,

542      G., DeReus, J., Janssen, S., Swafford, A.D., Orchanian, S.B., Sanders, J.G., Shorenstein, J.,

543      Holste, H., Petrus, S., Robbins-Pianka, A., Brislawn, C.J., Wang, M., Rideout, J.R., Bolyen,

544      E., Dillon, M., Caporaso, J.G., Dorrestein, P.C., Knight, R., 2018. QIITA: rapid, web-enabled

545      microbiome meta-analysis. Nat. Met. 15, 1–6. https://doi.org/10.1038/s41592-018-0141-

546      9

547   Hildebrand, F., 2021. Ultra-resolution metagenomics: when enough is not enough.

548      mSystems 6, e00881-21. https://doi.org/10.1128/msystems.00881-21

549   Humblot, C., Guyot, J.P., 2009. Pyrosequencing of tagged 16S rRNA gene amplicons for rapid

550      deciphering of the microbiomes of fermented foods such as pearl millet slurries. Appl.

551      Environ. Microbiol. 75, 4354–4361. https://doi.org/10.1128/aem.00451-09

552   Jagadeesan, B., Gerner-Smidt, P., Allard, M.W., Leuillet, S., Winkler, A., Xiao, Y., Chaffron, S.,

553      Vossen, J.V.D., Tang, S., Katase, M., McClure, P., Kimura, B., Chai, L.C., Chapman, J.,

554      Grant, K., 2019. The use of next generation sequencing for improving food safety:

555      Translation into practice. Food Microbiol. 79, 96–115.

556      https://doi.org/10.1016/j.fm.2018.11.005

557   Johnson, J.S., Spakowicz, D.J., Hong, B.-Y., Petersen, L.M., Demkowicz, P., Chen, L., Leopold,

558      S.R., Hanson, B.M., Agresta, H.O., Gerstein, M., Sodergren, E., Weinstock, G.M., 2019.

559      Evaluation of 16S rRNA gene sequencing for species and strain-level microbiome analysis.

560      Nat. Commun. 10, 5029. https://doi.org/10.1038/s41467-019-13036-1

561   Kovac, J., 2019. Precision food safety: a paradigm shift in detection and control of foodborne

562      pathogens. mSystems 4, e00164-19. https://doi.org/10.1128/msystems.00164-19

563     Kurtz, Z.D., Müller, C.L., Miraldi, E.R., Littman, D.R., Blaser, M.J., Bonneau, R.A., 2015. Sparse

564         and compositionally robust inference of microbial ecological networks. PloS Comput.

565         Biol. 11, e1004226. https://doi.org/10.1371/journal.pcbi.1004226

566     Lagkouvardos, I., Joseph, D., Kapfhammer, M., Giritli, S., Horn, M., Haller, D., Clavel, T., 2016.

567         IMNGS: A comprehensive open resource of processed 16S rRNA microbial profiles for

568         ecology and diversity studies. Sci. Rep. 6, 33721. https://doi.org/10.1038/srep33721

569     Lamprecht, A.-L., Garcia, L., Kuzak, M., Martinez, C., Arcila, R., Pico, E.M.D., Angel, V.D.D.,

570         Sandt, S. van de, Ison, J., Martinez, P.A., McQuilton, P., Valencia, A., Harrow, J.,

571         Psomopoulos, F., Gelpi, J.L., Hong, N.C., Goble, C., Capella-Gutierrez, S., 2020. Towards

572         FAIR principles for research software. Lect. Notes Comput. Sc. 3, 37–59.

573         https://doi.org/10.3233/ds-190026

574     McMurdie, P.J., Holmes, S., 2013. Phyloseq: an R package for reproducible interactive

575         analysis and graphics of microbiome census data. PloS One 8, e61217.

576         https://doi.org/10.1371/journal.pone.0061217

577     McMurdie, P.J., Holmes, S., 2015. Shiny-phyloseq: Web application for interactive

578         microbiome analysis with provenance tracking. Bioinformatics 31, 282–283.

579         https://doi.org/10.1093/bioinformatics/btu616

580     Meola, M., Rifa, E., Shani, N., Delbès, C., Berthoud, H., Chassard, C., 2019. DAIRYdb: a

581         manually curated reference database for improved taxonomy annotation of 16S rRNA

582         gene sequences from dairy products. BMC Genom. 20, 560.

583         https://doi.org/10.1186/s12864-019-5914-8

584     Mitchell, A.L., Almeida, A., Beracochea, M., Boland, M., Burgin, J., Cochrane, G., Crusoe,

585         M.R., Kale, V., Potter, S.C., Richardson, L.J., Sakharova, E., Scheremetjew, M.,

586         Korobeynikov, A., Shlemov, A., Kunyavskaya, O., Lapidus, A., Finn, R.D., 2019. MGnify: the

587 microbiome analysis resource in 2020. Nucl. Ac. Res. 48, D570–D578.

588 https://doi.org/10.1093/nar/gkz1035

589 Parente, E., Cocolin, L., De Filippis, F., Zotta, T., Ferrocino, I., O'sullivan, O., Neviani, E.,

590 Angelis, M.D., Cotter, P.D., Ercolini, D., 2016. FoodMicrobionet: A database for the

591 visualisation and exploration of food bacterial communities based on network analysis.

592 Int. J. Food Microbiol. 219, 28–37. https://doi.org/10.1016/j.ijfoodmicro.2015.12.001

593 Parente, E., De Filippis, F., Ercolini, D., Ricciardi, A., Zotta, T., 2019. Advancing integration of

594 data on food microbiome studies: FoodMicrobionet 3.1, a major upgrade of the

595 FoodMicrobionet database. Int. J. Food Microbiol. 305, 108249.

596 https://doi.org/10.1016/j.ijfoodmicro.2019.108249

597 Parente, E., Ricciardi, A., Zotta, T., 2020. The microbiota of dairy milk: a review. Int. Dairy J.

598 107, 104714. https://doi.org/10.1016/j.idairyj.2020.104714

599 Parente, E., Zotta, T., Ricciardi, A., 2021. Microbial association networks in cheese: a meta-

600 analysis. BiorXiv 2021.07.21.453196. https://doi.org/10.1101/2021.07.21.453196

601 Peschel, S., Müller, C.L., Mutius, E. von, Boulesteix, A.-L., Depner, M., 2021. NetCoMi:

602 network construction and comparison for microbiome data in R. Brief. Bioinform.

603 https://doi.org/10.1093/bib/bbaa290

604 Pollock, J., Glendinning, L., Wisedchanwet, T., Watson, M., 2018. The madness of

605 microbiome: attempting to find consensus "best practice" for 16S microbiome studies.

606 Appl. Environ. Microbiol. 84, 3225. https://doi.org/10.1128/aem.02627-17

607 R Core Team (2021). R: A language and environment for statistical computing. R Foundation

608 for Statistical Computing, Vienna, Austria. URL https://www.R-project.org/.

609    Ramakodi, M.P., 2022. Influence of 16S rRNA reference databases in amplicon-based

610        environmental microbiome research. Biotechnol Lett 44, 523–533.

611        https://doi.org/10.1007/s10529-022-03233-2

612    Roh, S.W., Kim, K.-H., Nam, Y.-D., Chang, H.-W., Park, E.-J., Bae, J.-W., 2010. Investigation of

613        archaeal and bacterial diversity in fermented seafood using barcoded pyrosequencing.

614        ISME J. 4, 1–16. https://doi.org/10.1038/ismej.2009.83

615    Shannon, P., Markiel, A., Ozier, O., Baliga, N.S., Wang, J.T., Ramage, D., Amin, N.,

616        Schwikowski, B., Ideker, T., 2003. Cytoscape: a software environment for integrated

617        models of biomolecular interaction networks. Genome Res. 13, 2498–2504.

618        https://doi.org/10.1101/gr.1239303

619    Smith, P.E., Waters, S.M., Expósito, R.G., Smidt, H., Carberry, C.A., McCabe, M.S., 2020.

620        Synthetic sequencing standards: a guide to database choice for rumen microbiota

621        amplicon sequencing analysis. Front. Microbiol. 11, 606825.

622        https://doi.org/10.3389/fmicb.2020.606825

623    Wang, L. G., Lam, T.T.Y., Xu, S ., Dai, Z ., Zhou, L ., Feng, T ., Guo, P ., Dunn, C.W., Jones, B.

624        R., Bradley, T ., Zhu, H ., Guan, Y., Jiang, Y., Yu, G. 2020. treeio: an R package for

625        phylogenetic tree input and output with richly annotated and associated data. Mol. Biol.

626        Evol. 37(2):599-603. doi: 10.1093/molbev/msz240

627    Werner, J.J., Koren, O., Hugenholtz, P., DeSantis, T.Z., Walters, W.A., Caporaso, J.G.,

628        Angenent, L.T., Knight, R., Ley, R.E., 2012. Impact of training sets on classification of high-

629        throughput bacterial 16s rRNA gene surveys. Isme J. 6, 94–103.

630        https://doi.org/10.1038/ismej.2011.82

631    Wilkinson, M.D., Dumontier, M., Aalbersberg, Ij.J., Appleton, G., Axton, M., Baak, A.,

632        Blomberg, N., Boiten, J.-W., Santos, L.B. da S., Bourne, P.E., Bouwman, J., Brookes, A.J.,

633      Clark, T., Crosas, M., Dillo, I., Dumon, O., Edmunds, S., Evelo, C.T., Finkers, R., Gonzalez-

634      Beltran, A., Gray, A.J.G., Groth, P., Goble, C., Grethe, J.S., Heringa, J., Hoen, P.A.C. 't,

635      Hooft, R., Kuhn, T., Kok, R., Kok, J., Lusher, S.J., Martone, M.E., Mons, A., Packer, A.L.,

636      Persson, B., Rocca-Serra, P., Roos, M., Schaik, R. van, Sansone, S.-A., Schultes, E.,

637      Sengstag, T., Slater, T., Strawn, G., Swertz, M.A., Thompson, M., Lei, J. van der, Mulligen,

638      E. van, Velterop, J., Waagmeester, A., Wittenburg, P., Wolstencroft, K., Zhao, J., Mons, B.,

639      2016. The FAIR Guiding Principles for scientific data management and stewardship. Sci.

640      Data 3, 160018. https://doi.org/10.1038/sdata.2016.18

641   Yap, M., Ercolini, D., Álvarez-Ordóñez, A., O'Toole, P.W., O'Sullivan, O., Cotter, P.D., 2021.

642      Next-generation food research: use of meta-omic approaches for characterizing

643      microbial communities along the food chain. Annu. Rev. Food Sci. Technol. 13, 1–24.

644      https://doi.org/10.1146/annurev-food-052720-010751

645   Yu, G. 2020. Using ggtree to visualize data on tree-like structures. Curr. Prot. Bioinf. 69:e96.

646      doi:10.1002/cpbi.96

647   Zotta, T., Ricciardi, A., Condelli, N., Parente, E., 2021. Metataxonomic and metagenomic

648      approaches for the study of undefined strain starters for cheese manufacture. Crit. Rev.

649      Food Sci. 1–15. https://doi.org/10.1080/10408398.2020.1870927

650