1  Comparative analysis of amplicon and metagenomic sequencing methods reveals key features

2  in the evolution of animal metaorganisms

3  Philipp Rausch[1,2,3,&,*], Malte Rühlemann[4,&], Britt Hermes[1,2,5], Shauni Doms[1,2], Tal Dagan[6], Katja

4  Dierking[7], Hanna Domin[8], Sebastian Fraune[8], Jakob von Frieling[9], Ute Henschel Humeida[10,11],

5  Femke-Anouska Heinsen[4], Marc Höppner[4], Martin Jahn[10,11], Cornelia Jaspers[11,12], Kohar Annie

6  B. Kissoyan[7], Daniela Langfeldt[6], Ateeqr Rehman[4], Thorsten B. H. Reusch[11,12], Thomas Röder[9],

7  Ruth A. Schmitz[6], Hinrich Schulenburg[7], Ryszard Soluch[6], Felix Sommer[4], Eva Stukenbrock[13,14],

8  Nancy Weiland-Bräuer[6], Philip Rosenstiel[4], Andre Franke[4], Thomas Bosch[8], John F. Baines[1,2,*]

9  Affiliations:

10  [1] Evolutionary Genomics, Max Planck Institute for Evolutionary Biology, Plön, Germany

11  [2] Institute for Experimental Medicine, Kiel University, Kiel, Germany

12  [3] Laboratory of Genomics and Molecular Biomedicine, Department of Biology University of
13  Copenhagen, Copenhagen Ø, Denmark

14  [4] Institute of Clinical Molecular Biology, Kiel University, Kiel, Germany

15  [5] Lübeck Institute of Experimental Dermatology, University of Lübeck, Germany

16  [6] Institute of General Microbiology, Kiel University, Kiel, Germany

17  [7] Department of Evolutionary Ecology and Genetics, Zoological Institute, Kiel University, Kiel,
18  Germany

19  [8] Zoological Institute, Kiel University, Kiel, Germany

20  [9] Molecular Physiology, Zoological Institute, Kiel University, Kiel, Germany

21  [10] Marine Ecology, Research Unit Marine Microbiology, GEOMAR Helmholtz Centre for Ocean
22  Research, Kiel, Germany

23  [11] Kiel University, Kiel, Germany

24  [12] Marine Ecology, GEOMAR Helmholtz Centre for Ocean Research, Kiel, Germany

25  [13] Environmental Genomics, Max Planck Institute for Evolutionary Biology, Plön, Germany

26  [14] Environmental Genomics, Botanical Institute, Kiel University, Kiel, Germany

27  & Authors contributed equally

28  * Corresponding authors: Philipp Rausch (philipp.rausch@bio.ku.dk), John F. Baines
29  (baines@evolbio.mpg.de)

30  Email addresses: John F. Baines- baines@evolbio.mpg.de, Thomas Bosch-
31  tbosch@zoologie.uni-kiel.de, Tal Dagan- tdagan@ifam.uni-kiel.de, Katja Dierking-
32  kdierking@zoologie.uni-kiel.de, Hanna Domin- hdomin@zoologie.uni-kiel.de, Shauni Doms-
33  doms@evolbio.mpg.de, Andre Franke- a.franke@mucosa.de, Sebastian Fraune-
34  sfraune@zoologie.uni-kiel.de, Jakob von Frieling- jfrieling@zoologie.uni-kiel.de, Femke-
35  Anouska Heinsen- f.heinsen@ikmb.uni-kiel.de, Ute Henschel Humeida-
36  uhentschel@geomar.de, Britt Marie Hermes- hermes@evolbio.mpg.de, Marc Höppner-
37  m.hoeppner@ikmb.uni-kiel.de, Martin Jahn- mjahn@geomar.de, Cornelia Jaspers-
38  cjaspers@geomar.de, Kohar Annie B. Kissoyan- kkissoyan@zoologie.uni-kiel.de, Daniela
39  Langfeldt- dlangfeldt@ifam.uni-kiel.de, Philipp Rausch- philipp.rausch@bio.ku.dk, Ateeqr
40  Rehman- a.rehmann@mucosa.de, Thorsten B. H. Reusch- treusch@geomar.de, Thomas
41  Röder- troeder@zoologie.uni-kiel.de, Philip Rosenstiel- p.rosenstiel@mucosa.de, Malte
42  Rühlemann- m.ruehlemann@ikmb.uni-kiel.de, Ruth A. Schmitz- rschmitz@ifam.uni-kiel.de,
43  Hinrich Schulenburg- hschulenburg@zoologie.uni-kiel.de, Ryszard Soluch- rsoluch@ifam.uni-
44  kiel.de, Felix Sommer- f.sommer@ikmb.uni-kiel.de, Eva Stukenbrock- estukenbrock@bot.uni-
45  kiel.de, Nancy Weiland-Bräuer- nweiland@ifam.uni-kiel.de

46

**Abstract**

**Background:** The interplay between hosts and their associated microbiome is now recognized as a fundamental basis of the ecology, evolution and development of both players. These interdependencies inspired a new view of multicellular organisms as "metaorganisms". The goal of the Collaborative Research Center "Origin and Function of Metaorganisms" is to understand why and how microbial communities form long-term associations with hosts from diverse taxonomic groups, ranging from sponges to humans in addition to plants.

**Methods:** In order to optimize the choice of analysis procedures, which may differ according to the host organism and question at hand, we systematically compared the two main technical approaches for profiling microbial communities, 16S rRNA gene amplicon- and metagenomic shotgun sequencing across our panel of ten host taxa. This includes two commonly used 16S rRNA gene regions and two amplification procedures, thus totaling five different microbial profiles per host sample.

**Conclusion:** While 16S rRNA gene-based analyses are subject to much skepticism, we demonstrate that many aspects of bacterial community characterization are consistent across methods and that metagenomic shotgun results are largely dependent on the employed pipeline. The resulting insight facilitates the selection of appropriate methods across a wide range of host taxa. Finally, by contrasting taxonomic and functional profiles and performing phylogenetic analysis, we provide important and novel insight into broad evolutionary patterns among metaorganisms, whereby the transition of animals from an aquatic to a terrestrial habitat marks a major event in the evolution of host-associated microbial composition.

**Keywords:** animal microbiome; evolution; phylosymbiosis; holobiont; metaorganism

**Background**

Dynamic host-microbe interactions have shaped the evolution of life. Virtually all plants and animals are colonized by an interdependent complex of microorganisms, and there is growing recognition that the biological processes of hosts and their associated microbial communities function in tandem, often as biological partners comprising a collective entity known as the metaorganism [1]. For instance, symbiotic bacteria contribute to host health and development in critical ways, ranging from nutrient metabolism to regulating whole life cycles [2] and in turn

78 benefit from habitats and resources the host provides. Moreover, it is well established that
79 perturbations of the microbiome likely play an important role in many host disease states [3].
80 However, researchers have yet to elucidate the mechanisms driving these interactions, as the
81 exact molecular and cellular processes are only poorly understood.

82 An integrated view on the metaorganism encompasses a cross-disciplinary approach
83 that addresses how and why microbial communities form long-term associations with their hosts.
84 Despite widespread agreement that the interdependencies of microbes and their hosts warrant
85 elucidation, there remains considerable incongruity between researchers regarding the best
86 methodologies to study host-microbe interactions. The development of standardized protocols
87 for characterizing and analyzing host-associated microbiomes across the breadth of the tree of
88 life are thus crucial to understand the evolution and function of metaorganisms without the
89 issues of technical inconsistencies or data quality.

90 Rapidly growing interest in microbiome research has been bolstered by the ability to
91 profile diverse microbial communities using next-generation sequencing (NGS). This culture-
92 free, high-throughput technology enables identification and comparison of entire microbial
93 communities [4]. Metagenomics typically encompasses two particular sequencing strategies:
94 amplicon sequencing, most often of the 16S rRNA gene as a phylogenetic marker, or shotgun
95 sequencing, which captures the complete breadth of DNA within a sample [4].

96 The use of the 16S ribosomal RNA gene as a phylogenetic marker has proven to be an
97 efficient and cost-effective strategy for microbiome analysis, and even allows for the imputation
98 of functional content based on taxon abundances [5]. However, PCR-based phylogenetic marker
99 protocols are vulnerable to biases through sample preparation and sequencing errors, in
100 particular the choice of which hypervariable regions of the 16S rRNA gene targeted seem to be
101 among the biggest factors underlying technical differences in microbiome composition [6-8].
102 Furthermore, 16S rRNA gene amplicon sequencing is typically limited to taxonomic classification
103 at the genus-level depending on the database and classifiers used [9], and provides only limited
104 functional information [5]. These well-recognized limitations of amplicon-based microbial
105 community analyses have raised concerns about the accuracy and reproducibility of 16S rRNA
106 phylogenetic marker studies and have led to an increased interest in developing more reliable
107 methods for amplicon library preparation and sequencing [8, 10].

108 Shotgun metagenomics, on the other hand, offers the advantage of species- and strain-
109 level classification of bacteria. Additionally, it allows researchers to examine the functional
110 relationships between hosts and bacteria by determining the functional content of samples

111  directly [9, 11], and enables the exploration of yet unknown microbial life that would otherwise

112  remain unclassifiable [12]. However, the relatively high costs of shotgun metagenomics and

113  more demanding bioinformatic requirements have precluded its use for microbiome analysis on

114  a wide scale [4, 9].

115  In this study, we set out to systematically compare experimental and analytical aspects of

116  the two main technical approaches for profiling microbial communities, 16S rRNA gene

117  amplicon- and shotgun sequencing, across a diverse array of host species studied in the

118  Collaborative Research Center 1182, "Origin and Function of Metaorganisms". The ten host

119  species range from basal aquatic metazoans [*Aplysina aerophoba* (sponge) and *Mnemiopsis*

120  *leidyi* (comb jelly)], to marine and limnic cnidarians (*Aurelia aurita, Nematostella vectensis*,

121  *Hydra vulgaris)*, standard vertebrate (*Mus musculus*) and invertebrate model organisms

122  (*Drosophila melanogaster*, *Caenorhabditis elegans*), to *Homo sapiens*, in addition to wheat

123  (*Triticum aestivum*) and a standardized mock community. This setup provides a breadth of

124  samples in terms of taxonomic composition and diversity. Conducting standardized data

125  generation procedures on these diverse samples on the one hand provides a unique and

126  powerful opportunity to systematically compare alternative methods, which display considerable

127  heterogeneity in performance. On the other hand, this information enables researchers working

128  on these or similar host species to choose the experimental (*e.g.* hypervariable region) or

129  analytical pipelines that best suit their needs, which will be a valuable resource to the greater

130  community of host-microbe researchers. Finally, we identified a number of interesting, broad

131  scale patterns contrasting the aquatic and terrestrial environment of metaorganisms, which also

132  reflect their evolutionary trajectories.

133

134  **Results**

135  Our panel of hosts includes ten species, for which five biological replicates each were included

136  (see Figure S1). The majority of hosts are metazoans, including the "gold sponge" (*Aplysina*

137  *aerophoba)*, moon jellyfish (*Aurelia aurita*), comb jellyfish (*Mnemiopsis leidyi*), starlet sea

138  anemone (*Nematostella vectensis)*, fresh-water polyp *Hydra vulgaris*, roundworm

139  (*Ceanorhabditis elegans),* fruit fly (*Drosophila melanogaster*), mouse (*Mus musculus*), human

140  (*Homo sapiens*), as well as the inclusion of wheat (*Triticum aestivum*), which can serve as an

141  outgroup to the metazoan taxa. *Drosophila melanogaster* was additionally sampled using two

142  different methods targeting feces and intestinal tissue. Nucleic acid extraction procedures were

143  conducted according to the needs of the individual host species (see Methods and

144    Supplementary Material), after which all DNA templates were subjected to a standard panel of
145    sequencing procedures. For 16S rRNA gene amplicon sequencing we used primers flanking two
146    commonly used variable regions, the V1V2 and V3V4 regions. Further, for each region we
147    compared a single-step fusion-primer PCR to a two-step procedure designed to improve the
148    accuracy of amplicon-based studies [8]. Finally, all samples were also subjected to shotgun
149    sequencing, such that five different sequence profiles were generated for each sample. While a
150    single classification pipeline was employed for all four 16S rRNA gene amplicon sequence
151    profiles, community composition based on shotgun data was initially evaluated using five
152    different classification methods (Kraken [13], MEGAN [14], MetaPhlan [15], MetaPhlan2 [16],
153    and SortmeRNA [17]; see Supplementary Material for comparative descriptions). However, due
154    to the advantage of simultaneously performing taxonomical and functional classification of
155    shotgun reads, as well as overall good performance (see analyses of mock community below),
156    MEGAN was used as a representative pipeline for most subsequent analyses.

157    **Performance of data processing and quality control:** All data generated from amplicons were
158    subject to the same stringent quality control pipeline including read-trimming, merging of forward
159    and reverse reads, quality filtering based on sequence quality and estimated errors, and chimera
160    removal (see Methods). The one step V1V2 amplicon data showed the highest rate of read-
161    survival (62.13 ± 23.90%, mean ± s.d.) followed by the corresponding two step method (mean=
162    49.85 ± 23.90%, mean ± s.d.), in large part due to the greater coverage of this comparatively
163    shorter amplicon (~312 bp). In contrast, 42.02 ± 16.41% and 36.88 ± 23.89% of the total reads
164    were included in downstream analysis for the one step and two step V3V4 data, respectively.
165    The longer V3V4 amplicon (~470 bp) was more affected by drops in quality at the end of the
166    reads, which decreases the overlap of forward and reverse reads and thus increases the
167    chances of sequencing errors (Figure S2, for final sample sizes see Table S1). Overall, aside
168    from chimera removal, each quality control step resulted in a comparatively greater loss of
169    V3V4- compared to V1V2 data. On the other hand, the V3V4 one step method yields the lowest
170    number of chimeras, suggesting a lower rate of chimera formation- and/or detection in this
171    approach (variable region- $F_{1,214}$=3.8881, $P$=0.0499, PCR- $F_{1,214}$=8.1751, $P$=0.0047, variable
172    region×PCR- $F_{1,214}$=6.4733, $P$=0.0117; Linear Mixed Model with organism as random factor).
173    Among all host taxa we observe the highest proportion of retained reads in the V1V2 one step
174    method and the lowest in the V3V4 two step method (Figure S2B; variable region-
175    $F_{1,215}$=74.9989, $P$<0.0001, PCR- $F_{1,215}$=21.0743, $P$<0.0001; Linear Mixed Model with organism
176    as random factor). After quality filtering and the identification of bacterial reads, an average of
177    0.46 Gb of shotgun reads per sample was achieved (range 0.03 to 2.1 Gb) (Figure S3A, for final

178    sample sizes see Table S1). To provide an initial assessment and comparison between the
179    amplicon and shotgun-based techniques, we plotted the discovered classifiable taxa and
180    functions for the entire pooled dataset. Although the methods differ distinctly, each method
181    shows a plateau in the number of discovered entities (see Figure S3C, S3D).

182    **Mock community:** The analysis of standardized mock communities is an important measure to
183    ensure general quality standards in microbial community analysis. In this study we employed a
184    commercially available mixture of eight bacterial- and two yeast species. Comparison among the
185    amplification procedures (one- and two step PCR), 16S rRNA gene regions (V1V2, V3V4) and
186    shotgun data reveals varying degrees of similarity to the expected microbial community
187    composition (Figure 1). One discrepancy is apparent due to the misclassification of
188    *Escherichia/Shigella,* whose close relationship make delineation at the genus level difficult
189    based on the V1V2 region are subsequently classified to *Enterobacteriaceae* (Figure 1A, Figure
190    S4). Classification of this bacterial group also differs according to shotgun pipeline employed,
191    due to different naming and taxonomic standards of the respective databases (*Escherichia*,
192    *Shigella*, *Enterobacteriaceae* refer to the *Escherichia/ Shigella* cluster) [18]. However, overall the
193    amplicon-based profiles show the closest matches to the expected community. The V1V2 one
194    step method and Kraken show the lowest degree of deviation between observed and expected
195    abundances of the focus taxa (Table 1, Figure S4). However, Kraken falsely detects a large
196    number of taxa not present in the mock communities. In addition, the relative abundances of
197    fungi in the mock community were relatively well predicted by MEGAN and Kraken, while
198    MetaPhlan2 failed to identify *Cryptococcus* and replaced it with several other taxa (see Figure
199    1).

200    Next, we evaluated alpha and beta diversity across the different technical and analytical
201    methods. Interestingly, most methods overestimate taxon richness but underestimate complexity
202    (as measured by the Shannon index) of the mock community, which could reflect biases arising
203    from grouping taxon abundances together (Figure 1, Figure S4, Figure S5, Table S2). Overall
204    the amplicon methods appear to more accurately reflect alpha diversity, although significant
205    differences are present with regard to the amplified region (species richness: variable region-
206    $F_{1,10}$=6.3657, $P$=0.0302; Shannon H: method- $F_{1,9}$=3.330, $P$=0.1014, variable region- $F_{1,9}$=6.110,
207    $P$=0.0354). With regard to beta diversity, the largest distance to the expected composition is
208    observed in SortmeRNA applied to shotgun sequencing of the mock community, while the
209    amplicon-based techniques, MEGAN, and MetaPhlan2 show the lowest distance (Figure 1D,
210    Figure S5, Table S3). Pairwise tests show almost no differences between the amplicon-based
211    techniques, while all shotgun based methods significantly differ from each other (Table S4).

212    Thus, in conclusion shotgun-based analysis pipelines yield a higher degree of variability/error
213    compared to the amplicon-based approaches based on a simple mock community. For
214    subsequent analyses we thus mainly focus on the amplicon-based data and MEGAN as a
215    representative shotgun-based pipeline, for which eukaryotic (*e.g.* fungal) sequences were not
216    included in the following analyses.

217    **Taxonomic diversity within and between hosts:** To evaluate the performance of our panel of
218    metagenomic methods over the range of complex host-associated communities in our
219    consortium, we next employed a panel of alpha- and beta diversity analyses to these samples,
220    which also provides an opportunity to infer broad patterns across animal taxa based on a
221    standardized methodology. Measures of alpha diversity display overall consistent values with
222    respect to host species, although many significant differences between technical methods are
223    present, mostly in a host-specific manner (Figure 2A-B). However, several host taxa display high
224    levels of consistency across methods including *A. aurita*, *C. elegans*, *D. melanogaster* and
225    *H. sapiens*, which show almost no significant differences between methods. Discrepancies and
226    individual recommendations for each host species are discussed in the Supplementary Material
227    (see Figures S6-S16). An intriguing observation is the tendency of aquatic hosts to display
228    higher alpha diversity values than those of terrestrial hosts, which is supported by average
229    differences between aquatic and terrestrial hosts and by relative consistent comparisons among
230    single host species as well (Figure 2C-D, Table S5). Finally, we also compared alpha diversity
231    estimates based on the other shotgun-based classifiers, which in most cases display greater
232    heterogeneity than among the 16S rRNA gene amplicon and MEGAN based estimates alone,
233    but still recover similar trends (Figure S17).

234    In order to investigate broad patterns of bacterial community similarity according to
235    metagenomic procedure and host species, we performed beta diversity analyses including all
236    host samples and each of their five different methodological profiles. This analysis reveals an
237    overall strong signal of host species, irrespective of the method used to generate community
238    profiles (Table 2, Figure 3). Pairwise comparisons between hosts are significant in all cases
239    except for samples derived from the V3V4 two step protocol, which did not consistently reach
240    significance after correction for multiple testing (Table S6). Further, complementary to the
241    observations made for alpha diversity, we also find strong signals of community differentiation
242    between the aquatic and terrestrial hosts (Table 2, Figure 3B and D). The separation between
243    these environments appears to be stronger based on amplicon data, whereas the separation
244    between hosts is stronger based on shotgun derived data (Table 2). Clustering of communities
245    based on host environment is consistent irrespective of the underlying shotgun analysis method,

246    although the topologies vary strongly (*e.g.* MetaPhlan2, see Figure S18). To further evaluate the

247    variability among biological replicates, we evaluated intra-group distances according to host

248    species, which reveals organisms with generally higher community variability (*i.e. C. elegans,*

249    *A. aurita, H. sapiens, H. vulgaris, T. aestivum*, and *M. leidyi*) than other host organisms in our

250    study (*N. vectensis, M. musculus, D. melanogaster, and A. aerophoba*; Figure S19A, C).

251    Interestingly, intra-group distances also significantly differ between the aquatic and terrestrial

252    environments, whereby aquatic organisms tend to display less variable communities than

253    terrestrial ones (Figure S19B, D). The low performance of *T. aestivum* in subsequent analyses

254    possibly originates from its commercial origin and low bacterial biomass relative to host material.

255         To identify individual drivers behind patterns of beta diversity, we performed indicator

256    species analysis [19] at the genus level with respect to method, host species, and environment.

257    Based on the amplicon data we identified 56 of 313 indicators to display consistent associations

258    across all four amplicon techniques, such as *Bacteroides, Barnesiella, Clostridium IV*, and

259    *Faecalibacterium* in *H. sapiens*, and *Helicobacter* and *Mucispirillum* in *M. musculus*, whereas

260    other associations were limited to *e.g.* only one variable region (Table S7, S8). However, the

261    overall pattern of host associations is largely consistent across methods (Figure S20). We also

262    identified numerous indicator genera for aquatic and terrestrial hosts (Table S9, S10). Indicator

263    analyses based on shotgun data reveals a smaller and less diverse set of host-specific

264    indicators, which however show many congruencies with the amplicon-based data.

265    **Functional diversity within and between hosts:** To examine the diversity (gene richness) of

266    metagenomic functions across host species we evaluated EggNOG [20] annotations (assembly-

267    based and MEGAN) to obtain a general functional spectrum (evolutionary genealogy of genes:

268    Non-supervised Orthologous Groups), in addition to annotations derived from a database

269    dedicated to functions interacting with carbohydrates (CAZY- Carbohydrate-Active enZYmes)

270    [21]. Overall the individual host communities differ drastically in gene richness (EggNOG genes

271    (MEGAN): $\chi^2$=52.202, *P*<2.10×10$^{-16}$; EggNOG genes (assembly): $\chi^2$=49.986, *P*<2.10×10$^{-16}$;

272    CAZY: $\chi^2$=48.815, *P*<2.10×10$^{-16}$; approximate Kruskal-Wallis test). Although the values also

273    differ considerably between methods, overall the functional repertoires are most diverse in the

274    vertebrate hosts, while only *H. vulgaris* and *A. aerophoba* as aquatic hosts carry a comparably

275    diverse functional repertoire (Figure 4A, Figure S21). Interestingly, in contrast to taxonomic

276    diversity we observe no difference in functional diversity between aquatic and terrestrial hosts.

277         Next we examined community differences (beta diversity) at the functional level, which

278    are overall more pronounced (average adj. $R^2$: 0.5084, Figure 4) than those based on taxonomic

279 (genus level) classification (shotgun adj. $R^2$: 0.4756; amplicon average adj. $R^2$: 0.4594, see
280 Table 2 and Table 3, Figure 3 and Figure 4, Figure S22). On the functional level aquatic and
281 terrestrial hosts are considerably less distinct than observed at the taxonomic level (taxonomic
282 shotgun adj. $R^2$=0.0766; taxonomic amplicon average adj. $R^2$=0.0690, functional shotgun
283 average adj. $R^2$=0.0441, see Table 2 and Table 3, Figure 4, S22). Variability of the functional
284 repertoires was lowest in *A. aerophoba*, *D. melanogaster* feces and *M. musculus* gut contents,
285 while *H. vulgaris*, *C. elegans*, and *D. melanogaster* gut samples displayed the highest intra-
286 group distances, which translates to a higher amount of functional heterogeneity between
287 replicates (Figure S23). This reflects in large part the patterns we observed in taxonomic
288 variability of those host-associated communities (Figure S19).

289 **Indicator functions:** To identify specific functions that are characteristic of individual hosts, we
290 applied indicator analysis to functional categories. General functions in EggNOG reveal several
291 interesting patterns, including CRISPR related genes in *A. aerophoba*, *H. sapiens*, and
292 *H. vulgaris*, suggesting a particular importance of viruses in these communities. *A. aerophoba*
293 possess a large set of characteristic genes involved in energy production and conversion, amino
294 acid transport and metabolism, replication, recombination and repair. *M. musculus* and others
295 appear to possess a large number of characteristic genes involved in carbohydrate transport
296 and metabolism, energy production and conversion, transcription and cell
297 wall/membrane/envelope biogenesis. *H. vulgaris* is characterized by a high number of genes
298 involved in transcription, inorganic ion transport, metabolism, signal transduction mechanisms
299 and cell wall/membrane/envelope biogenesis (Table S11-S13).

300 Analysis of carbohydrate-metabolizing functions based on CAZY [21] (Carbohydrate-
301 Active enZYmes) reveals the highest number of characteristic glycoside hydrolases (GH) in
302 *H. sapiens* and *M. musculus*, whereas polysaccharide lyases (PLs) for non-hydrolytic cleavage
303 of glycosidic bonds are present in *A. aerophoba* and *H. sapiens* (Table S14). Parts of the
304 cellulosome are only present in *A. aerophoba* and not in *M. musculus* or *H. sapiens*.
305 Interestingly, only the freshwater *H. vulgaris* carries characteristic auxiliary CAZYs involved in
306 lignin and chitin digestion, which may reflect dietary adaptations of the host.

307 **Performance of metagenome imputation from 16S rRNA gene amplicon data using**
308 **PICRUSt across metaorganisms:** Researchers often desire to obtain the insight gained from
309 functional metagenomic information despite being limited to 16S rRNA gene data, for which
310 imputation methods such as PICRUSt can be employed [5]. However, due to their dependence
311 on variable region and database coverage [5], these imputations must be viewed with caution.

312 Given our data set of both 16S amplicon- and shotgun metagenomic sequences, we
313 systematically evaluated the performance of PICRUSt predictions across hosts and amplicon
314 data type (V1V2, V3V4, one step/ two step protocol). Beginning with the mock community, the
315 V1V2 region displays lower performance for imputing functions compared to V3V4, as indicated
316 by a higher weighted Nearest Sequenced Taxon Index (NSTI) ($t$=17.812, $P$=1.119×10$^{-7}$, Figure
317 S24). High NSTI values imply low availability of genome representatives for the respective
318 sample, due to either large phylogenetic distance for each OTU to its closest sequenced
319 reference genome or a high frequency of poorly represented OTUs [5]. Comparing the
320 distribution of functional categories based on Clusters of Orthologous Groups (COG) [22]
321 between the different imputations (no cutoff applied) and the actual shotgun based repertoires
322 reveals considerable overlap (Figure S24). Exceptions include the functional category R
323 (general function prediction only), which is almost absent in the shotgun data, while the category
324 S (function unknown) is more abundant among the shotgun based functional data (Figure S24).

325 Next we evaluated functional imputations for the different host species and amplification
326 methods. We found no significant difference in average NSTI values or prediction success
327 (NSTI < 0.15) between amplification protocols or variable region. However, approximately a third
328 (31.8%) of the samples are lost due to incomplete imputation (NSTI > 0.15; Figure 5A). Notable
329 problematic host taxa are *A. aerophoba* and *H. vulgaris*, for which no sample remained below
330 the NSTI cutoff value. Other host taxa displayed clear differential performance with regard to the
331 variable region used, whereby *H. sapiens*, *N. vectensis* and *T. aestivum* were successfully
332 predicted based on V3V4, but not V1V2. However, when we employ Procrustes tests to
333 compare community functional profiles based on shotgun sequencing (single assembly,
334 MEGAN) and functional imputations at the COG-category level, we find a lower correspondence
335 of the V3V4-based imputations compared to those based on V1V2 (Figure 5B), while the
336 amplification methods displayed no significant difference. A similar pattern is observed when we
337 correlate community differences based on shotgun results and lower level (single functions)
338 COG annotations based on PICRUSt, although the difference is not significant ($F_{1,18}$=0.6172,
339 $P$=0.4423).

340 To investigate the similarities among methods in more detail, we merged shotgun and
341 PICRUSt based annotations at the level of COG categories. Principle coordinate analysis
342 reveals only small differences between imputations with regard to amplification method or
343 variable region (Figure 5C). However, large differences exist between the PICRUSt and shotgun
344 based functional repertoires, as well as between the shotgun techniques (MEGAN, single
345 assembly). Differences between the shotgun techniques were significant, but smaller than their

346 distance to the imputed functional spectra (Figure 5C, Table S15). Finally, we examined the
347 abundance of functional categories within single host taxa and the mock community, which
348 reveals a higher relative abundance of functions related to energy production and conversion
349 (C), replication, recombination and repair (L), and unknown functions (S) in the assembly-based
350 annotations compared to the other techniques, which might be an important driver of the
351 observed differences (Figure S24, S25).

352        Thus, in summary, the PICRUSt imputed functional repertoires significantly differ from
353 actual shotgun profiles. While variation in imputation success is largely dependent on the identity
354 of the particular host community, V3V4 appears to more often yield successful imputations.
355 However, when successful, V1V2-derived imputations display closer similarity to actual
356 functional profiles. Finally, the amplification method (one step, two step) appears to have no
357 significant effect on the quality of functional imputation. These data therefore support the notion
358 that metagenome imputations should be evaluated with care, as they depend on the underlying
359 variable region and sample source.

360 **Phylogenetic patterns in microbial community composition:** The term "phylosymbiosis"
361 refers to the phenomenon where the pattern of similarity among host-associated microbial
362 communities parallels the phylogeny of their hosts [23]. Highly divergent hosts with drastic
363 differences in physiology and life history might be expected to overwhelm the likelihood of
364 observing phylosymbiosis, which is typically observed within a given host clade [23]. However,
365 the factors driving differences in composition among our panel of hosts may also be expected to
366 vary in terms of the bacterial phylogenetic scale at which they are most readily observed [24].
367 Thus, we evaluated the degree to which bacterial community relationships (beta diversity) reflect
368 the underlying phylogeny of our hosts at a range of bacterial taxonomic ranks, spanning from the
369 genus to the phylum level.

370        In order to assess the general overlap between beta diversity and phylogenetic distance
371 of the host species, we performed Procrustes analysis [25]. These analyses reveal that the
372 strongest phylogenetic signal is observed when bacterial taxa are grouped at the order and/or
373 family level, whereby the one step protocols and the V3V4 region display greater correlations to
374 phylogenetic distance (Figure 6A). A similar pattern is observed for shotgun based community
375 profiles (*i.e.* MEGAN), although its fit increases again at the genus level. Measuring beta
376 diversity based on co-occurrence of bacterial taxa between hosts (Jaccard) displays a weaker
377 correspondence to host phylogeny than the abundance-based measure (Bray-Curtis) (Figure 6).

378    To assess the fit of individual host taxa, we examined the residuals of the correlation
379    between community composition and phylogenetic distance. This reveals a large variation in
380    correspondence among host taxa, with *M. musculus*, *M. leidyi*, *H. sapiens* and *D. melanogaster*
381    (feces) displaying the highest, while *H. vulgaris*, *C. elegans,* and *A. aerophoba* display the
382    lowest correspondence between their microbiome composition and phylogenetic position
383    (largest residuals Figure S26). Furthermore, terrestrial hosts display an overall better
384    correspondence between co-occurrences of bacterial genera and host relatedness (V1V2 one
385    step: *Z*=2.9578, *P*=0.0025), as do measurements based on V3V4 (one step: *Z*=2.7496,
386    *P*=0.0054; two step: *Z*=2.8097, *P*=0.0046; approximate Wilcoxon test).

387    Next, given the peak of correspondence between bacterial community composition and
388    host phylogeny observed at the order and/or family level, we set out to identify individual
389    community members whose abundances best correlate to host phylogenetic distance using
390    Moran's eigenvector method [26]. This reveals 41 bacterial families and 36 orders with
391    significant phylogenetic signal based on one or more amplicon data set, whereby 16 families
392    and 18 orders display repeated associations across methods (*e.g. Clostridia, Ruminococcaceae*,
393    *Helicobacteraceae*,        *Lachnospiraceae*,        *Coriobacteriaceae*,        *Erysipelotrichaceae*,
394    *Selenomonadales*, *Bacteroidales*, *Desulfovibrionales*; Table S16; Figure S27, S28). Analyzing
395    communities based on shotgun data on the other hand identifies 215 bacterial families and 97
396    orders associated with phylogenetic distances, whereby 69 and 27 display repeated
397    associations, respectively (Table S17; Figure S29, S30). The combined results of these
398    analyses identify several families and orders with strong and consistent phylogenetic
399    associations, in particular for the vertebrate hosts (*e.g. Bacteroidaceae/ Bacteroidales*,
400    *Bifidobacteriaceae/ Bifidobacteriales*, *Coriobacteriaceae/ Coriobacteriales*, *Desulfovibrionaceae/*
401    *Desulfovibrionales*,        *Erysipelotrichaceae/*        *Erysipelotrichales*,        *Porphyromonadaceae/*
402    *Bacteroidales*, *Ruminococcaceae/ Clostridiales*, *Selenomonadales*; see Table S16). Other
403    individual examples include bacteria related to *Helicobacteraceae/ Campylobacterales* in
404    *A. aurita*, which are observed in other marine cnidarians and may be involved in sulfur oxidation
405    [27]. *Alcanivoracaceae*, an alkane degrading bacterial group, is strongly associated to the
406    coastal cnidarian *N. vectensis*. This association might originate from adaptation to a polluted
407    coastal environment [28]. *Acidobacteria Gp6 and Gp9* specifically occur in *A. aerophoba* and are
408    commonly associated to the core microbial community of sponges [29].

409    **Phylogenetic patterns in functional community composition:** In order to contrast the
410    patterns observed at the taxonomic level to those based on function we used Procrustes
411    correlation to measure the overlap between phylogenetic distance and community distance

412 based on the panel of functional categories in our analyses. Interestingly, the two functional
413 categories displaying the greatest correspondence to host phylogeny are the CAZY and single
414 EggNOG based functions (Figure 6). The remainder of patterns between phylogeny and
415 bacterial functional spectra differed among the host species and functional categories (Figure
416 S26), *T. aestivum* and *D. melanogaster* (feces) display the lowest correspondence, while
417 *C. elegans, M. musculus and H. sapiens* display the best correspondence (lowest residuals,
418 Figure S26) between their functional repertoire and phylogenetic position. As observed for the
419 taxonomic analyses, terrestrial hosts again display a slightly better correlation than aquatic hosts
420 (smaller residuals), in particular for the co-abundance of EggNOG categories ($Z$=2.2116,
421 $P$=0.0267), CAZY ($Z$=2.0393, $P$=0.0414) and the co-occurrence of EggNOG categories
422 ($Z$=2.7377, $P$=0.0061) and genes ($Z$=3.3062, $P$=0.0007; approximate Wilcoxon test) among
423 hosts.

424       Finally, to reveal individual functions correlating to host phylogeny, we used the
425 aforementioned Moran's I eigenvector analyses with additional indicator analyses to narrow the
426 potential clade associations. Interestingly, most functions that correlate to a specific host
427 taxon/clade (1-3 taxa) are mainly restricted to vertebrate hosts or in combination with a
428 vertebrate host (Table S18-S21). This pattern is repeated across all functional annotations used
429 in this study. Examples include fucosyltransferases, fucosidases, polysaccharide binding
430 proteins, as well as hyaluronate, xanthan, and chondroitin lyases that stem from CAZY (see
431 Figure S31, Table S18). These functions are all related to glycan- and mucin degradation and
432 interaction, which mediate many intimate host-bacterial interactions and are also observed in
433 subsequent analyses based on general functional databases (EggNOG; Table S19, Table S20).
434 Many other phylogenetically correlated functions appear to be driven by the vertebrate hosts as
435 well, which likely reflects the high functional diversity within this group (see Figure 4 and Figure
436 S23). Only *LPXC* and *LPXK* (EggNOG), genes involved in the biosynthesis of the outer
437 membrane, are exclusively associated to the non-vertebrate hosts (LPXC: UDP-3-O-acyl-N-
438 acetylglucosamine deacetylase, LPXK: Tetraacyldisaccharide 4'-kinase), as is an oxidative
439 damage repair function (MSRA reductase) associated to *H. vulgaris* (Table S19, Figure S31).
440 EggNOG category Q (secondary metabolites biosynthesis, transport and catabolism) is also
441 characteristic of invertebrate hosts in addition to a small number of metabolic functions (*i.e.*
442 dehydrogenases, mono oxygenase, fatty acid hydroxylase; MEGAN based; Table S20, Figure
443 S31). More generally we observe a high number of genes of unknown function (S), carbohydrate
444 transport and metabolism (G), replication, recombination and repair (L), cell
445 wall/membrane/envelope biogenesis (M), and energy production and conversion (C) (Table S21

446 Figure S31). Finally, antibiotic resistance genes and virulence factors also show frequent
447 phylogenetic and host specific signals (Table S19, S20; Figure S31).

448

449 **Discussion**

450 Despite the great number of metagenomic studies published to date, which range in their focus
451 on technical, analytical or biological aspects, our study represents a unique contribution given its
452 breadth of different host samples analyzed with a panel of standardized methods. In particular,
453 the tradeoffs between 16S rRNA gene amplicon- versus shotgun sequencing concerning
454 amplification bias, functional information and both monetary and computational costs, warrant
455 careful consideration when designing research projects. While 16S rRNA gene amplicon-based
456 analyses are subject to considerable skepticism and criticism, we demonstrate that in many
457 aspects similar, if not superior characterization of bacterial communities is achieved by these
458 methods, although discrepancies associated with shotgun based data are largely dependent on
459 the analytic pipeline. We also show, however, that important insight can be gained through the
460 combination of taxonomic- and functional profiling, and that imputation-based functional profiles
461 significantly differ from actual profiles. Our findings thus provide a guide for selecting an
462 appropriate methodology for metagenomic analyses across a variety of metaorganisms. Finally,
463 these data provide novel insight into the broad scale evolution of host-associated bacterial
464 communities, which can be viewed as particularly reliable given the repeatability of observations
465 (*e.g.* differences between aquatic and terrestrial hosts, indicator taxa) across methods.

466 Given the concerns regarding the accuracy of 16S rRNA gene amplicon sequencing,
467 other studies such as that of Gohl *et al.* [8] performed systematic comparisons of different library
468 preparation methods, and found superior results for a two step amplification procedure. This
469 method offers the additional advantage that one panel of adapter/barcode sequences can be
470 combined with any number of different primers. Our first analyses were based on a standard
471 mock community including Gram positive and Gram negative bacteria from the Bacilli and
472 Gamma Proteobacteria (eight species), as well as two fungi, which did not support an
473 improvement of performance based on the two step protocol. However, a number of changes
474 were made to the Gohl *et al.* [8] protocol to adapt it to our lab procedures (*e.g.* larger reaction
475 volumes, polymerase, variable region, heterogeneity spacers) that may contribute to these
476 discrepancies, in addition to our different and diverse set of samples and other factors with
477 potential influence on the performance of amplicon sequencing [6-8, 30-32]. The complexity of
478 the mock community, *i.e.* the number of taxa, distribution, and phylogenetic breadth, may also

479   have an influence on the discovery of clear trends in amplification biases or detection limits for
480   certain taxonomic groups [33]. Thus, the even and phylogenetically shallow mock community in
481   our study may be less suited than the staggered and diverse mixtures used in other studies [8],
482   but still provides valuable information on repeatability, primer biases, and accuracy [33].
483   Nonetheless, when applied to our range of complex host-associated communities, we also found
484   that significant differences in most parameters were due to the variable region rather than
485   amplification method, and in many cases biological signals were either improved- or limited to
486   the one step protocol.

487   Additional sources of variation influencing the outcome of our 16S rRNA gene amplicon-
488   based community profiling are the bioinformatic pipelines we employed, starting from trimming
489   and merging to clustering and classification, which are stringent and incorporate more reliable
490   *de novo* clustering algorithms [34] as well as different classification databases [35].
491   Heterogeneity among the different amplicon approaches is however far smaller than the
492   observed heterogeneity between amplicon and shotgun methods, or within different shotgun
493   analyses, as observed in other benchmarking studies [31]. Differences between shotgun
494   approaches have been investigated in detail and also yield varying performances among
495   classifiers, but in general find a comparatively high performance of MEGAN based approaches
496   [9, 36, 37], which we also confirm in our study.

497   Given the limited number of studies that have compared imputed- and shotgun derived
498   functional repertoires [5, 38], our study also provides important additional insights. As imputation
499   by definition is data-dependent, the differential performance and prediction among hosts in our
500   study may in large part be explained by the amount of bacteria isolated, sequenced, and
501   deposited (16S rRNA or genome) from these hosts or their respective environments. This seems
502   to be most critical for the aquatic hosts. Furthermore, we observe a clear effect of variable region
503   on the prediction performance, which is most obvious based on the mock community. The
504   PICRUSt algorithm was developed and tested using primers targeting V3V4 16S rRNA, thus
505   optimization of the imputation algorithm might be biased towards this target over the V1V2
506   variable region. Although these performance differences, in particular the bias towards model
507   organisms compared to less characterized communities (*e.g.* hypersaline microbial mats), were
508   previously shown [5], our study provides additional, experimentally validated guidelines for a
509   number of novel host taxa.

510   Interestingly, the strongest correspondence between bacterial community similarity and
511   host genetic distance was detected at the bacterial order level for most of the employed

512  methods. This may on the one hand reflect the deep phylogenetic relationships between our

513  host taxa, such that turnover of bacterial taxa erodes phylosymbiosis over time [23, 24]. On the

514  other hand, some of the more striking observations made among our host taxa are the

515  differences between aquatic and terrestrial hosts, both at the level of alpha and beta diversity.

516  Based on a molecular clock for the 16S rRNA gene of roughly 1% divergence per 50 million

517  years [39], bacterial order level divergence corresponds well with the timing of animal

518  terrestrialization (425-500 MYA) [40, 41]. Although evolutionary rates can widely vary among

519  bacteria species [42], other studies of individual gut microbial lineages such as the *Enteroccoci*

520  indicate that animal terrestrialization was indeed a likely driver of diversification [43]. Specifically

521  the changing availability of carbohydrates in the host gut can be seen as a main driver of this

522  diversification, which is consistent with the association of CAZY-based functional repertoires

523  correlating to phylogenetic distance in our data set [23, 44].

524  In contrast to the patterns observed based on 16S rRNA gene amplicon-based profiles,

525  the differentiation of bacterial communities according to host habitat was less pronounced based

526  on functional genomic repertoires. This raises the possibility that the colonization of land by

527  ancient animals required the acquisition of new, land-adapted bacterial lineages to perform

528  some of the same ancestral functions. The overall observation of increased beta diversity among

529  terrestrial- compared to aquatic hosts (Figure S19) could in part reflect differential acquisition

530  among host lineages after colonizing land, although dispersal in the aquatic environment may on

531  the other hand act as a greater homogenizing factor among aquatic hosts. The stronger

532  correspondence between bacterial community- and host phylogenetic distance among terrestrial

533  hosts is also generally consistent with this hypothesis. However, the higher alpha diversity and

534  the slightly lower correspondence with the phylogenetic patterns in aquatic hosts may also

535  indicate a higher influence of environmental bacteria or a lack of physiological control over

536  bacterial communities.

537  Bacterial taxa and functions involved in carbohydrate utilization were among the most

538  notable associations to individual hosts, groups of hosts, and/or host phylogenetic relationships.

539  Taxa such as *Bacteroidales*, *Ruminococcaceae/ Ruminococcales*, and *Clostridia* associated to

540  humans and/or mice include members known for a mucosal lifestyle, and these hosts also

541  display the most diverse and abundant repertoire of carbohydrate active enzymes (particularly

542  glycosylhydrolases) in their microbiome. Other examples include sialidases, esterases, and

543  fucosyltransferases, as well as different extracellular structures that appear to be specific to

544  aquatic hosts, indicating differences in mucus and glycan composition according to this host

545  environment. Glycan structures provide a direct link between the microbial community and the

546   host via attachment, nutrition, and communication [45, 46], and the composition of mucin and

547   glycan structures themselves show strong evolutionary patterns and are distinct among

548   taxonomic groups [44]. Thus, a high diversity of glycan structures within and between hosts may

549   determine the specific sets carbohydrate facilitating enzymes of the respective microbial

550   communities.

551   In addition to the bacterial carbohydrate hydrolases that digest surrounding host and

552   dietary carbohydrates, we also identified a number of glycosyltransferases associated with

553   capsular polysaccharide synthesis (Table S19, Table S20). This type of glycosylation is an

554   important facilitator for host association and survival [47] and plays a crucial role in infections

555   [48]. The capsule prevents opsonization and phagocytosis through the host immune system and

556   gives the bacterium the ability to modulate its interaction with the host environment [47, 49]. This

557   type of manipulation is performed by mutualists and pathogens alike [47, 50] via molecular

558   mimicry and tolerogenic immune modulation [51, 52]. Bacterial glycan products like

559   polysaccharide A (PSA) may also have direct benefits for the host, as it can interfere with the

560   host immune system by increasing immunologic tolerance, or inhibit the binding of other

561   microbes (*e.g. Helicobacter hepaticus* [53]). Thus, capsular and excreted glycan structures are

562   important for the successful colonization and persistence in different environments [54, 55] and

563   host organisms [47, 55].

564

565   **Conclusions**

566   In summary, the systematic comparison of five different metagenomic sequencing

567   methods applied to ten different holobiont yielded a number of novel technical and biological

568   insights. Although important exceptions will exist, we demonstrate that broad scale biological

569   patterns are largely consistent across these varying methods. While the richer information

570   provided by shotgun sequencing is clearly desirable and is likely to surpass amplicon-based

571   profiling techniques in the foreseeable future, technical variability among analytical pipelines

572   currently surpasses that observed between different amplicon methods. As many aspects of

573   differential performance in our study are host-specific (more detailed description of individual

574   hosts can be found in the Supplementary Material), future development and benchmarking

575   analyses would also benefit from a including a range of different host/environmental samples.

576

577   **Methods**

578 **DNA extraction and 16S rRNA gene amplicon sequencing**: Protocols for each host type are
579 described in the Supplementary Material (see also Figure S18-S28). Each library (16S rRNA
580 gene amplicon, shotgun) included at least one mock community sample based on the
581 ZymoBIOMICS™ Microbial Community DNA Standard (Lot.: ZRC187324, ZRC187325)
582 consisting of 8 bacterial species (*Pseudomonas aeruginosa* (10.4%), *Escherichia coli* (9.0%),
583 *Salmonella enterica* (11.8%), *Lactobacillus fermentum* (10.3%), *Enterococcus faecalis* (14.1%),
584 *Staphylococcus aureus* (14.6%), *Listeria monocytogenes* (13.2%), *Bacillus subtilis* (13.2%)) and
585 two fungi (*Saccharomyces cerevisiae* (1.6%), *Cryptococcus neoformans* (1.8%)).

586 The 16S rRNA gene was amplified using uniquely barcoded primers flanking the V1 and
587 V2 hypervariable regions (27F-338R) and V3V4 hypervariable regions (515F-806R) with fused
588 MiSeq adapters and heterogeneity spacers in a 25 µl PCR [32]. For the traditional one step PCR
589 protocol we used 4 µl of each forward and reverse primer (0.28 µM), 0.5 µl dNTPs (200 µM
590 each), 0.25 µl Phusion Hot Start II High-Fidelity DNA Polymerase (0.5 Us), 5 µl of HF buffer
591 (Thermo Fisher Scientific, Inc., Waltham, MA, USA) and 1 µl of undiluted DNA. PCRs were
592 conducted with the following cycling conditions (98°C-30s, 30×[98°C-9s, 55°C-60s, 72°C-90s],
593 72°C-10 min) and checked on a 1.5 % agarose gel. Using a modified version of the recently
594 published two step PCR protocol by Gohl *et al.* 2016, we employed for the first round of
595 amplification fusion primers consisting of the 16S rRNA gene primers (V1V2, V3V4) and a part
596 of the Illumina Nextera adapter with the following cycling conditions in a 25 µl PCR reaction
597 (98°C-30s, 25×[98°C-10s, 55°C-30s, 72°C-60s], 72°C-10 min) [8]. Following the PCR was
598 diluted 1:10 and 5µl of the solution were used in an additional reaction of 10 µl (98°C-30s,
599 10×[98°C-9s, 55°C-30s, 72°C-60s], 72°C-10 min) utilizing the Nextera adapter overhangs to
600 ligate the Illumina adapter sequence and individual MIDs to the amplicons following the
601 manufacturer's instructions. The PCR protocol we used 1 µl of each forward and reverse primer
602 (5 µM), 0.3 µl dNTPs (10 µM), 0.2 µl Phusion Hot Start II High-Fidelity DNA Polymerase (2 U/µl),
603 2 µl of 5×HF buffer (Thermo Fisher Scientific, Inc., Waltham, MA, USA) and 5 µl of the diluted
604 PCR product. The concentration of the amplicons was estimated using a Gel Doc™ XR+
605 System coupled with Image Lab™ Software (BioRad, Hercules, CA USA) with 3 µl of
606 O'GeneRulerTM 100 bp Plus DNA Ladder (Thermo Fisher Scientific, Inc., Waltham, MA, USA)
607 as the internal standard for band intensity measurement. The samples of individual gels were
608 pooled into approximately equimolar subpools as indicated by band intensity and measured with
609 the Qubit dsDNA br Assay Kit (Life Technologies GmbH, Darmstadt, Germany). Sub pools were
610 mixed in an equimolar fashion and stored at -20°C until sequencing.

611   Library preparation for shotgun sequencing was performed using the NexteraXT kit
612   (Illumina) for fragmentation and multiplexing of input DNA following the manufacturer's
613   instructions. Amplicon sequencing was performed on the Illumina MiSeq platform with v3
614   chemistry (2×300 cycle kit), while shotgun sequencing was performed via 2×150bp Mid Output
615   Kit at the IKMB Sequencing Center (CAU Kiel, Germany).

616   **Amplicon analysis:** The respective V1V2 and V3V4 PCR primer sequences were removed
617   from the sequencing data using *cutadapt* (v.1.8.3) [56]. Sequence data in FastQ format was
618   quality trimmed using *sickle* (v.1.33) in paired-end mode with default settings and removing
619   sequences dropping below 100bp after trimming [57]. Forward and reverse read were merged
620   into a single amplicon read using VSEARCH allowing fragments with a length of 280-350 bp for
621   V1V2 and 350-500 bp for V3V4 amplicons [58]. Sequence data was quality controlled using
622   fastq_quality_filter (FastX Toolkit) retaining sequences with no more than 5% of per-base quality
623   values below 30 and subsequently with VSEARCH discarding sequences with more than 1
624   expected errors [58, 59]. Reference guided chimera removal was performed using the gold.fa
625   reference in VSEARCH (v2.4.3). The UTAX algorithm was used for a fast classification of the
626   sequence data in order to remove sequences not assigned to the domains Bacteria or Archaea
627   and exclude amplicon fragments from Chloroplasts [60]. Notably, only a total of 15 sequences
628   were assigned to the domain Archaea, all found in two samples of human feces, accounting for
629   less than 0.1% of the clean reads in theses samples. The entire cleaned sequence data was
630   concatenated into a single file, dereplicated and processed with VSEARCH for OTU picking
631   using the UCLUST algorithm [61] using a 97% similarity threshold. OTUs were again checked
632   for chimeric sequences, now using the *de novo* implementation of the UCHIME algorithm in
633   VSEARCH [58, 61, 62]. All clean sequence data of the samples were mapped back to the
634   cleaned OTU sequences using VSEARCH. OTU sequences and clean sequences mapping to
635   the OTUs were taxonomically annotated using the RDP classifier algorithm with the RDP training
636   set 14 [63, 64]. Sequence data were normalized by selecting 10,000 random sequences per
637   sample. Taxon-by-sample abundance tables were created for all taxonomic levels from Phylum
638   to Genus, as well as for OTUs.

639   **PICRUSt functional imputations:** Species level OTUs (97% similarity threshold) were further
640   classified using the GreenGenes (August 2013) database [65] via RDP classifier as
641   implemented in mothur (v1.39.5) and merged with the abundances into a biome file which was
642   uploaded to the Galaxy PICRUSt v1.1.1 pipeline (http://galaxy.morganlangille.com/) to derive
643   functional imputations (COG predicitions) [5]. To achieve accurate functional predictions

644 samples with NSTI ≤ 0.15 (weighted Nearest Sequenced Taxon Index) were pruned from the
645 data set, as recommended by the developers.

646 **Shotgun sequencing:** Raw demultiplexed sequences were trimmed via Trimmomatic (v0.36)
647 for low quality regions with a minimum length of 50 bp as well as for adaptor and remaining MID
648 sequences [66]. After trimming reads were mapped to host specific genome databases and *ΦX*
649 with additional retention databases containing all fully sequenced bacterial and metagenomic
650 genomes (05-09-2015) via DeconSeq (v0.4.3) [67]. Single and paired sequences were repaired
651 using the BBTools (v37.28) repair function [68]. Combined sequences were searched against
652 the non-redundant NCBI database (28-07-2017) via DIAMOND [69] with (evalue cutoff 0.001,
653 v0.8.28) and MEGAN [14] classifying hits by functions (EGGNOG-Oct2016) and taxa (May2017)
654 (v6.6.1). MetaPhlan [15] (v1.7.7) and MetaPhlan2 [16] (v2.2.0) was used for taxonomic
655 classification. Forward and reverse reads were mapped to the SIVLA non-redundant database
656 (v123) via SortmeRNA [17, 70] (2.1b) and classified via RDP classifier and the RPD 16 database
657 as implemented in mothur [71]. Kraken (v0.10.5-beta) database was constructed on complete
658 and dusted genome sequences of all archaea (+scaffolds), bacteria, fungi (+scaffolds), protozoa
659 (+scaffolds), viruses and full sequences of plasmids and plastids [13] (database 21-08-2017),
660 which was used to classify raw reads as well as assembled contigs, which were used throughout
661 the manuscript. For assemblies of single samples we used metaSPADES [72] (v3.9.1) using
662 paired reads in addition to unpaired reads left from the previous steps. PROKKA (v1.12) was
663 used for gene calling and initial genome annotation [73] using the metagenome option with
664 additional identifying rRNAs and snRNA via barnap, ARAGORN [74], and Infernal [75]. ORFs
665 were further annotated via EggNOG annotation via HMMER models implemented in the eggnog-
666 mapper (v0.12.7) [20, 76], CAZY database via dbCAN (v5, 07/24/2016) and HMMER3 [21, 77].
667 Gene abundances were derived from mapping the all reads back to the predicted ORF via
668 bowtie2 (v2.2.6) [78] and calculated TPM (transcripts per kilobase million) via SamTools (v1.5)
669 [79].

670 18S rRNA genes were obtained from NCBI GeneBank and aligned via ClustalW (v1.4)
671 [80] for host tree construction, which includes *A. aerophoba* (gi:51095211, AY5917991),
672 *M. leidyi* (gi:14517703, AF2937001), *H. vulgaris* (gi:761889987, JN5940542), *A. aurita*
673 (gi:14700050, AY0392081), *N. vectensis* (gi:13897746, AF2543821), *T. aestivum* (gi:15982656,
674 AY0490401), *M. musculus* (gi:374088232, NR_0032783), *H. sapiens* (gi:36162, X032051),
675 *D. melanogaster* (gi:939630477, NR_1335591), and *C.elegans* (gi:30525807, AY2681171).
676 Phylogenetic distance was calculated via DNADIST (v3.5c) [81] and a maximum likelihood tree
677 was constructed via FastTree v2.1 CAT+Γ model [82]. Accuracy was improved via increased

678    minimum evolution rounds for initial tree search [-spr 4], more exhaustive tree search [-mlacc 2],

679    and a slow initial tree search [-slownni].

680    **Statistical analysis:** Statistical analyses were carried via R [83] (v3.4.3). Alpha diversity indices

681    (richness, Shannon-Weaver index) and beta diversity metrics based on the shared presence

682    (Jaccard distance)- or abundance (Bray-Curtis distance) of taxa were calculated in the *vegan*

683    package [84] and ordinated via Principal Coordinate Analysis (PCoA, avoiding negative

684    eigenvalues), or via non-metric multidimensional scaling (NMDS) using a maximum of 10000

685    random starts to obtain a minimally stressed configuration in three dimensions. Clusters were fit

686    via an iterative process (10'000 permutations) tested for separation by direct gradient analysis

687    via distance based Redundancy analyses and permutative ANOVA (10'000 permutations) [85,

688    86]. Univariate analyses were carried out with approximate Wilcoxon/Kruskal tests as

689    implemented in *coin* [87] (10'000 permutations). Procrustes tests were used to relate pairwise

690    community distances based on either different data sources such as functional repertoires or

691    taxonomic composition, as well as phylogenetic distances [25, 88]. Moran's I eigenvector

692    technique was employed to correlate bacterial community members and their functions to

693    phylogenetic divergence, as implemented in *ape* (10'000 permutations) [26, 89] . Indicator

694    species analysis, employing the generalized indicator value (*IndVal.g*), was used to assess the

695    predictive value of a taxon for each respective host phenotype/category as implemented in

696    *indicspecies* [19]. Linear mixed models, as implemented in *nlme* were used to compare the

697    influence of amplification method or variable region without the influence of the organism of

698    origin [90]. We employed the Hommel- and Benjamini-Yekutieli adjustment of *P*-values when

699    advised [91, 92].

700

**Declarations**

**Ethics approval and consent to participate (Human samples):** Study participants were randomly recruited from inhabitants of Schleswig-Holstein (Germany) which were recruited for the PopGen cohort. Five individuals from the PopGen biobank (Schleswig-Holstein, Germany) were randomly selected among the healthy and unmedicated individuals and included in the study without corresponding meta-information. Study participants collected fecal samples at home in standard fecal tubes and shipped them immediately at room temperature or brought them to the collection center (within 24 h). Samples were stored at –80°C until processing. Human feces (N=4) were sampled and extracted following the procedures as described in Wang *et al.* 2016 [93]. A biopsy sample of the sigmoid colon was taken from a healthy control individual without macro- or microscopical inflammation (N=1) and DNA was extracted as described in Rausch *et al.* 2011 [94]. Investigators were blinded to sample identities and written, informed consent was obtained from all study participants before the study. All protocols were approved by the Ethics Committee of the Medical Faculty of Kiel and by the data protection officer of the University Hospital Schleswig-Holstein in adherence with the Declaration of Helsinki Principles.

**Ethics approval for animal and plant samples:** Wild derived, hybrid mice were sacrificed according to the German animal welfare law and Federation of European Laboratory Animal Science Associations guidelines. Hybrid breeding stocks of wild derived *M. m. musculus* × *M. m. domesticus* hybrids captured in 2008 are kept at the Max Planck Institute Plön (11th lab generation). The approval for mouse husbandry and experiment was obtained from the local veterinary office "Veterinäramt Kreis Plön" (Permit: 1401-144/PLÖ-004697). All sampling, including invertebrate and plant samples, was performed in concordance with the German animal welfare law and Federation of European Laboratory Animal Science Associations guidelines. Further details for each host type are provided in the Supplementary Material.

**Consent for publication:** Not applicable.

**Availability of data and material:** Sequence- and meta-data are accessible under the study identifier PRJEB30924 ("https://www.ebi.ac.uk/ena"). Remaining DNA from non-human samples can be made available upon request. All human samples and information on their corresponding phenotypes have to be obtained from the PopGen Biobank Kiel (Schleswig-Holstein, Germany) through a Material Data Access Form. Information about the Material Data Access Form and how to apply can be found at: "https://www.uksh.de/p2n/Information+for+Researchers.html".

734    **Competing Interests:** The authors declare no competing interests.

737    **Author contributions:** PRa, PRo, AF, TB, and JFB conceived and designed research. PRa and

738    MR performed data analyses. PRa, MR, BH, SD, and JFB interpreted results and wrote the

739    manuscript. PRa, MR, TD, KD, HD, SD, SF, JF, UHH, FAH, BH, MH, MJ, CJ, KABK, DL, AR,

740    TBHR, TR, RAS, HS, RS, FS, ES, NWB, PRo, AF, TB, and JFB generated and interpreted host-

741    specific data and gave intellectual input. All authors read and approved the final manuscript.

745

746  **Figure legends:**

747  **Figure 1:** Average community composition of bacteria (A) and fungi (B) in the mock community

748  samples sequenced via metagenomic shotgun- and 16S rRNA gene amplicon techniques

749  (amplicon: V1V2, V3V4, one step, two step; shotgun: MEGAN based classification (short reads),

750  MetaPhlan (short reads), MetaPhlan2 (short reads), Kraken based classification (contigs),

751  SortmeRNA (short reads)). (C) Bacterial genus-level alpha diversity estimates in comparison to

752  the expected community value. (D) Principle coordinate analysis of the Bray-Curtis distance

753  between methods and the expected community. Ellipses represent standard deviations of points

754  within the respective groups. Sample sizes for the different approaches are $N_{shotgun}=4$, $N_{V1V2\text{-}one\ step}=3$,

755  $N_{V1V2\text{-}two\ step}=3$, $N_{V3V4\text{-}one\ step}=3$, and $N_{V1V2\text{-}two\ step}=3$.

756  **Figure 2:** Comparison of bacterial genus richness (A) and Shannon H (B) based on 16S rRNA

757  gene amplicon and shotgun derived genus profiles based on MEGAN highlighting the

758  differences between variable regions, amplification methods, and metagenomic classifier, as

759  well as between the different host organisms. Colors show significance of amplification methods

760  (A, C) or pairwise comparisons of methods (B, D) based on pairwise *t*-tests with Hommel *P*-

761  value adjustment (A, B), and approximate Wilcoxon test for the comparison between

762  environmental categories (C, D). Mean values are shown in grey symbols in plots A and B.

763  Sample sizes are indicated below the samples.

764  **Figure 3:** Non-metric Multidimensional Scaling of Bray-Curtis distances based on genus profiles

765  derived from the different 16S rRNA gene amplicon methods (V1V2/ V3V4, one step/ two step)

766  and shotgun derived genus profiles highlighting (A) host differences and (B) differences between

767  host environments (terrestrial/aquatic; see Table 2). Non-metric Multidimensional Scaling of

768  Jaccard distances based on genus profiles derived from the different 16S rRNA gene amplicon

769  methods and shotgun derived genus profiles highlighting (C) host taxon differences and (D)

770  differences between host environments (terrestrial/aquatic; see Table 2). Both panels show a

771  separation based on host organisms and environments and not by method. Large symbols

772  indicate the centroid of the respective host groups and vertical lines help to determine their

773  position in space. Samples sizes are equal to Figure 2 (see also Table S1).

774  **Figure 4:** Multivariate correlation (Procrustes analyses) of phylogenetic distance among host

775  organisms and community distances based on 16S rRNA gene amplicon- or shotgun derived

776  community profiles at different taxonomic cutoffs, from Phylum to Genus and species level OTUs

777  in the amplicon based profiles. Similar results are shown for the correspondence between

778  functional composition based distances derived from imputed COGs and COG categories

779    imputed from PICRUSt, and EggNOG derived genes and COG categories, as well as CAZY. All

780    correlations are significant at $P \leq 0.05$ (10'000 permutations). Large symbols indicate the

781    centroid of the respective host groups and vertical lines help to determine their position in space.

782    **Figure 5:** (A) Differences in Nearest Sequenced Taxon Index (imputation success) between

783    variable regions (average: $Z=0.3869$, $P=0.7017$, approximate Wilcoxon test; probability: odds

784    ratio=1.5941, $P=0.1402$, Fisher test) and amplification method ($Z=0.0667$, $P=0.9472$,

785    approximate Wilcoxon test; probability: odds ratio=1.5511, $P=0.1436$, Fisher test). (B)

786    Procrustes correlation of imputed and shotgun based COG categories among different

787    techniques, with significantly higher correspondence between imputed and measured functional

788    profiles in the V1V2 compared to the V3V4 region ($F_{1,18}=7.8537$, $P=0.0118$, ANOVA). (C) Non-

789    metric Multidimensional Scaling displays Bray-Curtis distances based on functional category

790    abundances (COG categories) derived from PICRUSt (V1V2/ V3V4, one step/ two step) and

791    shotgun based approaches (MEGAN, single assembly). Ellipses represent standard deviations

792    of points within the respective groups.

793    **Figure 6:** Functional diversities were derived from the number and abundances of MEGAN

794    based EggNOG annotations. Functional richness between (**A**) host organisms and (**B**) host

795    environmental groups based is displayed, as well as functional differences between hosts (**C**)

796    and environmental groups (**D**). Non-metric Multidimensional Scaling is based on Bray-Curtis

797    distances on the differences in functional composition between the host organisms is displayed

798    (**C**, **D**; see Table 3). Large symbols indicate the centroid of the respective groups. Functional

799    variation of communities based on pairwise Bray-Curtis distances within host organism groups

800    and environmental groups. Samples sizes for the host taxa is N=5, except for D. melanogaster

801    gut tissue (N=10; see Table S1).

802 **Tables:**

803 **Table 1:** Differences between expected and observed genus abundances in the mock communities ($N_{shotgun}$=4, $N_{amplicon}$=3) via a one-

804 sample $t$-test (two-sided) of relative abundances ($P$-values are adjusted via Hommel procedure).

| Members mock community | shotgun | | | | | amplicon | | | |
|---|---|---|---|---|---|---|---|---|---|
| | MEGAN | Kraken | MetaPhlan | MetaPhlan2 | SortmeRNA | V1V2 one step | V3V4 one step | V1V2 two step | V3V4 two step |
| *Staphylococcus* | 0.00002 | 0.14039 | 0.07916 | 0.07010 | $1.1097 \times 10^{-6}$ | 0.52446 | 0.09200 | 0.03994 | 0.21564 |
| *Listeria* | 0.00395 | 0.06065 | 0.02306 | 0.06043 | $1.4751 \times 10^{-6}$ | 0.34964 | 0.53267 | 0.03003 | 0.00545 |
| *Bacillus* | 0.00006 | 0.09558 | 0.02219 | 0.03638 | $3.9824 \times 10^{-7}$ | 0.21420 | 0.02818 | 0.29671 | 0.30589 |
| *Pseudomonas* | 0.13668 | 0.40989 | 0.62649 | 0.46933 | $2.7877 \times 10^{-7}$ | 0.36721 | 0.05776 | 0.38147 | 0.59037 |
| *Escherichia/Shigella** | NA | NA | NA | NA | $9.9378 \times 10^{-10}$ | 0.00462 | 0.45612 | 0.00237 | 0.59037 |
| *Shigella** | $4.6372 \times 10^{-10}$ | NA | $8.0806 \times 10^{-8}$ | NA | NA | NA | NA | NA | NA |
| *Escherichia** | 0.00001 | 0.00882 | 0.00710 | 0.28178 | NA | NA | NA | NA | NA |
| *Enterobacteriaceae** | NA | NA | NA | NA | NA | 0.87898 | 0.00004 | 0.19274 | 0.00055 |
| *Salmonella* | $3.8092 \times 10^{-6}$ | 0.08772 | 0.02203 | 0.03361 | $4.9348 \times 10^{-7}$ | 0.34964 | 0.05838 | 0.09712 | 0.08851 |
| *Lactobacillus* | 0.00297 | 0.09704 | 0.05384 | 0.04043 | $1.4751 \times 10^{-6}$ | 0.87898 | 0.53267 | 0.38147 | 0.59037 |
| *Enterococcus* | 0.00012 | 0.18719 | 0.00353 | 0.07277 | $6.3719 \times 10^{-7}$ | 0.04816 | 0.03746 | 0.01159 | 0.00954 |

805 * Escherichia/Shigella relatives counted as equivalent

806

807 **Table 2:** Taxonomic distance based PERMANOVA results for differences in community composition (genus level) between host species

808 and host environments based on shared abundance (Bray-Curtis) and shared presence (Jaccard), based on whole genome shotgun and

809 different amplicon strategies ($P$-values are adjusted via Hommel procedure).

| Distance | Factor | Data | Classifier | $DF$ | $F$ | $P$ | $P_{Hommel}$ | $R^2$ | adj. $R^2$ |
|---|---|---|---|---|---|---|---|---|---|
| Bray-Curtis | organism | shotgun | MEGAN | 10,49 | 6.3517 | 0.0001 | 0.0001 | 0.5645 | 0.4756 |
| | | amplicon | V1V2-one step | 10,43 | 7.1026 | 0.0001 | 0.0001 | 0.6229 | 0.5352 |
| | | | V1V2-two step | 10,42 | 4.2297 | 0.0001 | 0.0001 | 0.5018 | 0.3831 |
| | | | V3V4-one step | 10,43 | 7.8964 | 0.0001 | 0.0001 | 0.6474 | 0.5654 |
| | | | V3V4-two step | 10,41 | 3.7917 | 0.0001 | 0.0001 | 0.4805 | 0.3538 |

| Distance | Factor | | Data | | DF | F | P | $P_{\text{Hommel}}$ | $R^2$ | adj. $R^2$ |
|---|---|---|---|---|---|---|---|---|---|---|
| | environment | shotgun | MEGAN | | 1,58 | 5.8958 | 0.0001 | 0.0004 | 0.0923 | 0.0766 |
| | | amplicon | V1V2-one step | | 1,52 | 6.1588 | 0.0001 | 0.0001 | 0.1059 | 0.0887 |
| | | | V1V2-two step | | 1,51 | 4.6185 | 0.0001 | 0.0001 | 0.0830 | 0.0651 |
| | | | V3V4-one step | | 1,52 | 5.4975 | 0.0001 | 0.0001 | 0.0956 | 0.0782 |
| | | | V3V4-two step | | 1,50 | 3.3349 | 0.0001 | 0.0001 | 0.0625 | 0.0438 |
| Jaccard | organism | shotgun | MEGAN | | 10,49 | 4.7458 | 0.0001 | 0.0001 | 0.4920 | 0.3883 |
| | | amplicon | V1V2-one step | | 10,43 | 3.6867 | 0.0001 | 0.0001 | 0.4616 | 0.3364 |
| | | | V1V2-two step | | 10,42 | 2.9760 | 0.0001 | 0.0001 | 0.4147 | 0.2754 |
| | | | V3V4-one step | | 10,43 | 4.0248 | 0.0001 | 0.0001 | 0.4835 | 0.3633 |
| | | | V3V4-two step | | 10,41 | 2.9343 | 0.0001 | 0.0001 | 0.4171 | 0.2750 |
| | environment | shotgun | MEGAN | | 1,58 | 4.3872 | 0.0001 | 0.0004 | 0.0703 | 0.0543 |
| | | amplicon | V1V2-one step | | 1,52 | 3.8714 | 0.0001 | 0.0001 | 0.0693 | 0.0514 |
| | | | V1V2-two step | | 1,51 | 3.6541 | 0.0001 | 0.0001 | 0.0669 | 0.0486 |
| | | | V3V4-one step | | 1,52 | 4.3213 | 0.0001 | 0.0001 | 0.0767 | 0.0590 |
| | | | V3V4-two step | | 1,50 | 3.6646 | 0.0001 | 0.0001 | 0.0683 | 0.0497 |

**Table 3:** Functional distance based PERMANOVA results for differences in general functional community composition (EggNOG) and carbohydrate active enzymes (CAZY) between host species and host environments based on shared abundance (Bray-Curtis) and shared presence (Jaccard) of functions (*P*-values are adjusted via Hommel procedure).

| Distance | Factor | Data | DF | F | P | $P_{\text{Hommel}}$ | $R^2$ | adj. $R^2$ |
|---|---|---|---|---|---|---|---|---|
| Bray-Curtis | organism | CAZY | 10,47 | 7.3323 | 0.0001 | 0.0001 | 0.6094 | 0.5263 |
| | | EggNOG categories | 10,49 | 5.6088 | 0.0001 | 0.0001 | 0.5337 | 0.4386 |
| | | EggNOG gene+description | 10,49 | 4.4454 | 0.0001 | 0.0001 | 0.4757 | 0.3687 |
| | | EggNOG (MEGAN categories) | 10,49 | 12.2594 | 0.0001 | 0.0001 | 0.7144 | 0.6562 |
| | | EggNOG (MEGAN gene) | 10,49 | 8.2788 | 0.0001 | 0.0001 | 0.6282 | 0.5523 |
| | environment | CAZY | 1,56 | 5.4257 | 0.0001 | 0.0007 | 0.0883 | 0.0721 |
| | | EggNOG categories | 1,58 | 2.5429 | 0.0195 | 0.0195 | 0.0420 | 0.0255 |
| | | EggNOG gene+description | 1,58 | 3.0662 | 0.0001 | 0.0007 | 0.0502 | 0.0338 |
| | | EggNOG (MEGAN categories) | 1,58 | 3.7703 | 0.0015 | 0.0030 | 0.0610 | 0.0448 |
| | | EggNOG (MEGAN gene) | 1,58 | 3.7271 | 0.0002 | 0.0012 | 0.0604 | 0.0442 |
| Jaccard | organism | CAZY | 10,47 | 3.9098 | 0.0001 | 0.0001 | 0.4541 | 0.3380 |
| | | EggNOG categories | 10,49 | 3.7179 | 0.0001 | 0.0001 | 0.4314 | 0.3154 |

810

811

812

813

|  |  |  |  |  |  |  |  |
|---|---|---|---|---|---|---|---|
|  | EggNOG gene+description | 10,49 | 2.5275 | 0.0001 | 0.0001 | 0.3403 | 0.2057 |
|  | EggNOG (MEGAN categories) | 10,49 | 7.7781 | 0.0001 | 0.0001 | 0.6135 | 0.5346 |
|  | EggNOG (MEGAN gene) | 10,49 | 5.4989 | 0.0001 | 0.0001 | 0.5288 | 0.4326 |
| environment | CAZY | 1,56 | 2.5866 | 0.0003 | 0.0021 | 0.0442 | 0.0271 |
|  | EggNOG categories | 1,58 | 1.4180 | 0.1442 | 0.1442 | 0.0239 | 0.0070 |
|  | EggNOG gene+description | 1,58 | 1.9535 | 0.0004 | 0.0024 | 0.0326 | 0.0159 |
|  | EggNOG (MEGAN categories) | 1,58 | 3.0425 | 0.0460 | 0.0920 | 0.0498 | 0.0335 |
|  | EggNOG (MEGAN gene) | 1,58 | 3.1222 | 0.0001 | 0.0009 | 0.0511 | 0.0347 |

814

815

**References:**

1.  Bosch TCG, McFall-Ngai MJ: **Metaorganisms as the new frontier**. *Zoology* 2011, **114**(4):185-190.
2.  McFall-Ngai M, Hadfield MG, Bosch TCG, Carey HV, Domazet-Lošo T, Douglas AE, Dubilier N, Eberl G, Fukami T, Gilbert SF *et al*: **Animals in a bacterial world, a new imperative for the life sciences**. *Proceedings of the National Academy of Sciences* 2013, **110**(9):3229-3236.
3.  Carding S, Verbeke K, Vipond DT, Corfe BM, Owen LJ: **Dysbiosis of the gut microbiota in disease**. *Microbial Ecology in Health and Disease* 2015, **26**(1):26191.
4.  Morgan XC, Huttenhower C: **Chapter 12: Human Microbiome Analysis**. *PLoS Comput Biol* 2012, **8**(12):e1002808.
5.  Langille MGI, Zaneveld J, Caporaso JG, McDonald D, Knights D, Reyes JA, Clemente JC, Burkepile DE, Vega Thurber RL, Knight R *et al*: **Predictive functional profiling of microbial communities using 16S rRNA marker gene sequences**. *Nat Biotech* 2013, **31**(9):814-821.
6.  Hiergeist A, Glasner J, Reischl U, Gessner A: **Analyses of Intestinal Microbiota: Culture versus Sequencing**. *ILAR journal / National Research Council, Institute of Laboratory Animal Resources* 2015, **56**(2):228-240.
7.  Tremblay J, Singh K, Fern A, Kirton ES, He S, Woyke T, Lee J, Chen F, Dangl JL, Tringe SG: **Primer and platform effects on 16S rRNA tag sequencing**. *Frontiers in Microbiology* 2015, **6**:771.
8.  Gohl DM, Vangay P, Garbe J, MacLean A, Hauge A, Becker A, Gould TJ, Clayton JB, Johnson TJ, Hunter R *et al*: **Systematic improvement of amplicon marker gene methods for increased accuracy in microbiome studies**. *Nat Biotech* 2016, **34**(9):942-949.
9.  Walsh AM, Crispie F, O'Sullivan O, Finnegan L, Claesson MJ, Cotter PD: **Species classifier choice is a key consideration when analysing low-complexity food microbiome data**. *Microbiome* 2018, **6**(1):50.
10. Faith JJ, Guruge JL, Charbonneau M, Subramanian S, Seedorf H, Goodman AL, Clemente JC, Knight R, Heath AC, Leibel RL *et al*: **The Long-Term Stability of the Human Gut Microbiota**. *Science* 2013, **341**(6141):1237439.
11. Jovel J, Patterson J, Wang W, Hotte N, O'Keefe S, Mitchel T, Perry T, Kao D, Mason AL, Madsen KL *et al*: **Characterization of the Gut Microbiome Using 16S or Shotgun Metagenomics**. *Frontiers in Microbiology* 2016, **7**(459):459.
12. Rinke C, Schwientek P, Sczyrba A, Ivanova NN, Anderson IJ, Cheng J-F, Darling A, Malfatti S, Swan BK, Gies EA *et al*: **Insights into the phylogeny and coding potential of microbial dark matter**. *Nature* 2013, **499**(7459):431-437.
13. Wood D, Salzberg S: **Kraken: ultrafast metagenomic sequence classification using exact alignments**. *Genome Biology* 2014, **15**(3):R46.
14. Huson D, Auch A, Qi J, Schuster S: **MEGAN analysis of metagenomic data**. *Genome Res* 2007, **17**(3):377 - 386.
15. Segata N, Waldron L, Ballarini A, Narasimhan V, Jousson O, Huttenhower C: **Metagenomic microbial community profiling using unique clade-specific marker genes**. *Nat Meth* 2012, **9**(8):811-814.
16. Truong DT, Franzosa EA, Tickle TL, Scholz M, Weingart G, Pasolli E, Tett A, Huttenhower C, Segata N: **MetaPhlAn2 for enhanced metagenomic taxonomic profiling**. *Nat Meth* 2015, **12**(10):902-903.
17. Kopylova E, Noe L, Touzet H: **SortMeRNA: fast and accurate filtering of ribosomal RNAs in metatranscriptomic data**. *Bioinformatics* 2012, **28**(24):3211-3217.

865    18.   Hong Nhung P, Ohkusu K, Mishima N, Noda M, Monir Shah M, Sun X, Hayashi M, Ezaki
866           T: **Phylogeny and species identification of the family Enterobacteriaceae based on**
867           **dnaJ sequences**. *Diagnostic Microbiology and Infectious Disease* 2007, **58**(2):153-161.
868    19.   De Cáceres M, Legendre P, Moretti M: **Improving indicator species analysis by**
869           **combining groups of sites**. *Oikos* 2010, **119**(10):1674-1684.
870    20.   Huerta-Cepas J, Szklarczyk D, Forslund K, Cook H, Heller D, Walter MC, Rattei T,
871           Mende DR, Sunagawa S, Kuhn M *et al*: **eggNOG 4.5: a hierarchical orthology**
872           **framework with improved functional annotations for eukaryotic, prokaryotic and**
873           **viral sequences**. *Nucleic Acids Research* 2016, **44**(1):286-293.
874    21.   Cantarel BL: **The Carbohydrate-Active EnZymes database (CAZy): an expert**
875           **resource for glycogenomics**. *Nucleic Acids Res* 2009, **37**(Database issue):233-238.
876    22.   Tatusov RL, Natale DA, Garkavtsev IV, Tatusova TA, Shankavaram UT, Rao BS,
877           Kiryutin B, Galperin MY, Fedorova ND, Koonin EV: **The COG database: new**
878           **developments in phylogenetic classification of proteins from complete genomes**.
879           *Nucleic Acids Res* 2001, **29**(1):22-28.
880    23.   Brooks AW, Kohl KD, Brucker RM, van Opstal EJ, Bordenstein SR: **Phylosymbiosis:**
881           **Relationships and Functional Effects of Microbial Communities across Host**
882           **Evolutionary History**. *PLOS Biology* 2016, **14**(11):e2000225.
883    24.   Groussin M, Mazel F, Sanders JG, Smillie CS, Lavergne S, Thuiller W, Alm EJ:
884           **Unraveling the processes shaping mammalian gut microbiomes over evolutionary**
885           **time**. *Nature Communications* 2017, **8**:14319.
886    25.   Peres-Neto P, Jackson D: **How well do multivariate data sets match? The**
887           **advantages of a Procrustean superimposition approach over the Mantel test**.
888           *Oecologia* 2001, **129**(2):169-178.
889    26.   Gittleman JL, Kot M: **Adaptation: Statistics and a Null Model for Estimating**
890           **Phylogenetic Effects**. *Systematic Zoology* 1990, **39**(3):227-241.
891    27.   Murray AE, Rack FR, Zook R, Williams MJM, Higham ML, Broe M, Kaufmann RS, Daly
892           M: **Microbiome Composition and Diversity of the Ice-Dwelling Sea Anemone,**
893           **Edwardsiella andrillae**. *Integrative and Comparative Biology* 2016, **56**(4):542-555.
894    28.   Schneiker S, dos Santos VAPM, Bartels D, Bekel T, Brecht M, Buhrmester J, Chernikova
895           TN, Denaro R, Ferrer M, Gertler C *et al*: **Genome sequence of the ubiquitous**
896           **hydrocarbon-degrading marine bacterium Alcanivorax borkumensis**. *Nature*
897           *Biotechnology* 2006, **24**(8):997.
898    29.   Hentschel U, Piel J, Degnan SM, Taylor MW: **Genomic insights into the marine**
899           **sponge microbiome**. *Nat Rev Micro* 2012, **10**(9):641-654.
900    30.   Wu JY, Jiang XT, Jiang YX, Lu SY, Zou F, Zhou HW: **Effects of polymerase, template**
901           **dilution and cycle number on PCR based 16 S rRNA diversity analysis using the**
902           **deep sequencing method**. *BMC Microbiol* 2010, **10**:255.
903    31.   D'Amore R, Ijaz UZ, Schirmer M, Kenny JG, Gregory R, Darby AC, Shakya M, Podar M,
904           Quince C, Hall N: **A comprehensive benchmarking study of protocols and**
905           **sequencing platforms for 16S rRNA community profiling**. *BMC Genomics* 2016,
906           **17**(1):55.
907    32.   Fadrosh D, Ma B, Gajer P, Sengamalay N, Ott S, Brotman R, Ravel J: **An improved**
908           **dual-indexing approach for multiplexed 16S rRNA gene sequencing on the Illumina**
909           **MiSeq platform**. *Microbiome* 2014, **2**(1):6.
910    33.   Highlander S: **Mock Community Analysis**. In: *Encyclopedia of Metagenomics.* Edited
911           by Nelson EK. New York, NY: Springer New York; 2013: 1-7.
912    34.   Westcott SL, Schloss PD: **De novo clustering methods outperform reference-based**
913           **methods for assigning 16S rRNA gene sequences to operational taxonomic units**.
914           *PeerJ* 2015, **3**:e1487.

915   35.   Werner JJ, Koren O, Hugenholtz P, DeSantis TZ, Walters WA, Caporaso JG, Angenent
916         LT, Knight R, Ley RE: **Impact of training sets on classification of high-throughput**
917         **bacterial 16s rRNA gene surveys**. *ISME J* 2012, **6**(1):94-103.
918   36.   Sczyrba A, Hofmann P, Belmann P, Koslicki D, Janssen S, Dröge J, Gregor I, Majda S,
919         Fiedler J, Dahms E *et al*: **Critical Assessment of Metagenome Interpretation—a**
920         **benchmark of metagenomics software**. *Nature Methods* 2017, **14**(11):1063–1071.
921   37.   Lindgreen S, Adair KL, Gardner PP: **An evaluation of the accuracy and speed of**
922         **metagenome analysis tools**. *Sci Rep* 2016, **6**:19233.
923   38.   Xu Z, Malmer D, Langille MGI, Way SF, Knight R: **Which is more important for**
924         **classifying microbial communities: who's there or what they can do?** *ISME J* 2014,
925         **8**(12):2357-2359.
926   39.   Ochman H, Elwyn S, Moran NA: **Calibrating bacterial evolution**. *Proceedings of the*
927         *National Academy of Sciences of the United States of America* 1999, **96**(22):12638-
928         12643.
929   40.   Benton MJ: **The origins of modern biodiversity on land**. *Philosophical Transactions of*
930         *the Royal Society B: Biological Sciences* 2010, **365**(1558):3667-3679.
931   41.   Rota-Stabelli O, Daley Allison C, Pisani D: **Molecular Timetrees Reveal a Cambrian**
932         **Colonization of Land and a New Scenario for Ecdysozoan Evolution**. *Current*
933         *Biology* 2013, **23**(5):392-398.
934   42.   Kuo CH, Ochman H: **Inferring clocks when lacking rocks: the variable rates of**
935         **molecular evolution in bacteria**. *Biol Direct* 2009, **4**:35.
936   43.   Lebreton F, Manson AL, Saavedra JT, Straub TJ, Earl AM, Gilmore MS: **Tracing the**
937         **Enterococci from Paleozoic Origins to the Hospital**. *Cell* 2017, **169**(5):849-861.e813.
938   44.   Bishop JR, Gagneux P: **Evolution of carbohydrate antigens—microbial forces**
939         **shaping host glycomes?** *Glycobiology* 2007, **17**(5):23R-34R.
940   45.   Pickard JM, Maurice CF, Kinnebrew MA, Abt MC, Schenten D, Golovkina TV, Bogatyrev
941         SR, Ismagilov RF, Pamer EG, Turnbaugh PJ *et al*: **Rapid fucosylation of intestinal**
942         **epithelium sustains host-commensal symbiosis in sickness**. *Nature* 2014,
943         **514**(7524):638-641.
944   46.   Schwartzman JA, Koch E, Heath-Heckman EAC, Zhou L, Kremer N, McFall-Ngai MJ,
945         Ruby EG: **The chemistry of negotiation: Rhythmic, glycan-driven acidification in a**
946         **symbiotic conversation**. *Proceedings of the National Academy of Sciences* 2015,
947         **112**(2):566-571.
948   47.   Martens EC, Chiang HC, Gordon Jl: **Mucosal Glycan Foraging Enhances Fitness and**
949         **Transmission of a Saccharolytic Human Gut Bacterial Symbiont**. *Cell Host &amp;*
950         *Microbe* 2008, **4**(5):447-457.
951   48.   Boulnois GJ, Roberts IS: **Genetics of capsular polysaccharide production in**
952         **bacteria**. *Current Topics in Microbiology and Immunology* 1990, **150**:1-18.
953   49.   Meng D, Newburg DS, Young C, Baker A, Tonkonogy SL, Sartor RB, Walker WA,
954         Nanthakumar NN: **Bacterial symbionts induce a FUT2-dependent fucosylated niche**
955         **on colonic epithelium via ERK and JNK signaling**. *American Journal of Physiology -*
956         *Gastrointestinal and Liver Physiology* 2007, **293**(4):780-787.
957   50.   Mahdavi J, Pirinccioglu N, Oldfield NJ, Carlsohn E, Stoof J, Aslam A, Self T, Cawthraw
958         SA, Petrovska L, Colborne N *et al*: **A novel O-linked glycan modulates**
959         **Campylobacter jejuni major outer membrane protein-mediated adhesion to human**
960         **histo-blood group antigens and chicken colonization**. *Open Biology* 2014,
961         **4**(1):130202.
962   51.   Severi E, Hood D, Thomas G: **Sialic acid utilization by bacterial pathogens**.
963         *Microbiology* 2007, **153**(Pt 9):2817 - 2822.
964   52.   Coyne MJ, Chatzidaki-Livanis M, Paoletti LC, Comstock LE: **Role of glycan synthesis**
965         **in colonization of the mammalian gut by the bacterial symbiont Bacteroides**
966         **fragilis**. *Proceedings of the National Academy of Sciences* 2008, **105**(35):13099-13104.

967   53.   Mazmanian SK, Round JL, Kasper DL: **A microbial symbiosis factor prevents**
968         **intestinal inflammatory disease**. *Nature* 2008, **453**(7195):620-625.
969   54.   Tounkang S, Premkumar D, Gustavo S, Nathalie B, Yann B, Patricia C, Florence L,
970         Olivier N, Brigitte G, Anne L *et al*: **Capsular glucan and intracellular glycogen of**
971         **Mycobacterium tuberculosis: biosynthesis and impact on the persistence in mice**.
972         *Molecular Microbiology* 2008, **70**(3):762-774.
973   55.   Roberts IS: **The biochemistry and genetics of capsular polysaccharide production**
974         **in bacteria**. *Annu Rev Microbiol* 1996, **50**(1):285-315.
975   56.   Martin M: **Cutadapt removes adapter sequences from high-throughput sequencing**
976         **reads**. *2011* 2011, **17**(1).
977   57.   Joshi N, Fass J: **Sickle: A sliding-window, adaptive, quality-based trimming tool for**
978         **FastQ files**. In., 1.33 edn. https://github.com/najoshi/sickle; 2011.
979   58.   Rognes T, Flouri T, Nichols B, Quince C, Mahé F: **VSEARCH: a versatile open source**
980         **tool for metagenomics**. *PeerJ* 2016, **4**:e2584.
981   59.   Gordon A, Hannon G: **Fastx-toolkit. FASTQ/A short-reads pre-processing tools**. In:
982         *Unpublished Available online at: http://hannonlab cshl edu/fastx_toolkit.* 2010.
983   60.   Edgar RC: **UTAX algorithm**. In.; 2015.
984   61.   Edgar RC: **Search and clustering orders of magnitude faster than BLAST**.
985         *Bioinformatics* 2010, **26**(19):2460-2461.
986   62.   Edgar RC, Haas BJ, Clemente JC, Quince C, Knight R: **UCHIME improves sensitivity**
987         **and speed of chimera detection**. *Bioinformatics* 2011, **27**(16):2194-2200.
988   63.   Wang Q, Garrity GM, Tiedje JM, Cole JR: **Naive Bayesian Classifier for Rapid**
989         **Assignment of rRNA Sequences into the New Bacterial Taxonomy**. *Applied and*
990         *environmental microbiology* 2007, **73**(16):5261-5267.
991   64.   Cole JR, Chai B, Marsh TL, Farris RJ, Wang Q, Kulam SA, Chandra S, McGarrell DM,
992         Schmidt TM, Garrity GM *et al*: **The Ribosomal Database Project (RDP-II): previewing**
993         **a new autoaligner that allows regular updates and the new prokaryotic taxonomy**.
994         *Nucl Acids Res* 2003, **31**(1):442-443.
995   65.   McDonald D, Price MN, Goodrich J, Nawrocki EP, DeSantis TZ, Probst A, Andersen GL,
996         Knight R, Hugenholtz P: **An improved Greengenes taxonomy with explicit ranks for**
997         **ecological and evolutionary analyses of bacteria and archaea**. *ISME J* 2012,
998         **6**(3):610-618.
999   66.   Bolger AM, Lohse M, Usadel B: **Trimmomatic: a flexible trimmer for Illumina**
1000        **sequence data**. *Bioinformatics* 2014, **30**(15):2114-2120.
1001  67.   Schmieder R, Edwards R: **Fast Identification and Removal of Sequence**
1002        **Contamination from Genomic and Metagenomic Datasets**. *PLoS One* 2011,
1003        **6**(3):e17288.
1004  68.   Bushnell B, Rood J: **BBTools bioinformatics tools, including BBMap**. In: *URL*
1005        *http://sourceforge net/projects/bbmap.* 37.28 edn; 2017.
1006  69.   Buchfink B, Xie C, Huson DH: **Fast and sensitive protein alignment using DIAMOND**.
1007        *Nat Meth* 2015, **12**(1):59-60.
1008  70.   Pruesse E, Quast C, Knittel K, Fuchs BM, Ludwig W, Peplies J, Glockner FO: **SILVA: a**
1009        **comprehensive online resource for quality checked and aligned ribosomal RNA**
1010        **sequence data compatible with ARB**. *Nucl Acids Res* 2007, **35**(21):7188-7196.
1011  71.   Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski RA,
1012        Oakley BB, Parks DH, Robinson CJ *et al*: **Introducing mothur: Open Source,**
1013        **Platform-independent, Community-supported Software for Describing and**
1014        **Comparing Microbial Communities**. *Applied and environmental microbiology* 2009,
1015        **75**(23):7537-7541.
1016  72.   Nurk S, Meleshko D, Korobeynikov A, Pevzner P: **metaSPAdes: a new versatile de**
1017        **novo metagenomics assembler**. *arXiv preprint arXiv:160403071* 2016.

1018 73.   Seemann T: **Prokka: rapid prokaryotic genome annotation**. *Bioinformatics* 2014,
1019        **30**(14):2068-2069.
1020 74.   Laslett D, Canback B: **ARAGORN, a program to detect tRNA genes and tmRNA**
1021        **genes in nucleotide sequences**. *Nucleic Acids Research* 2004, **32**(1):11-16.
1022 75.   Kolbe DL, Eddy SR: **Fast filtering for RNA homology search**. *Bioinformatics* 2011,
1023        **27**(22):3102-3109.
1024 76.   Huerta-Cepas J, Forslund K, Coelho LP, Szklarczyk D, Jensen LJ, von Mering C, Bork P:
1025        **Fast Genome-Wide Functional Annotation through Orthology Assignment by**
1026        **eggNOG-Mapper**. *Molecular Biology and Evolution* 2017, **34**(8):2115-2122.
1027 77.   Yin Y, Mao X, Yang J, Chen X, Mao F, Xu Y: **dbCAN: a web resource for automated**
1028        **carbohydrate-active enzyme annotation**. *Nucleic Acids Research* 2012, **40**(1):445-
1029        451.
1030 78.   Langmead B, Salzberg SL: **Fast gapped-read alignment with Bowtie 2**. *Nature*
1031        *methods* 2012, **9**(4):357-359.
1032 79.   Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G,
1033        Durbin R, Genome Project Data Processing S: **The Sequence Alignment/Map format**
1034        **and SAMtools**. *Bioinformatics* 2009, **25**(16):2078-2079.
1035 80.   Thompson JD, Higgins DG, Gibson TJ: **CLUSTAL W: improving the sensitivity of**
1036        **progressive multiple sequence alignment through sequence weighting, position-**
1037        **specific gap penalties and weight matrix choice**. *Nucleic Acids Research* 1994,
1038        **22**(22):4673-4680.
1039 81.   Felsenstein J: **DNADIST -- Program to compute distance matrix from nucleotide**
1040        **sequences**. In., 3.5c edn; 1993.
1041 82.   Price MN, Dehal PS, Arkin AP: **FastTree 2 – Approximately Maximum-Likelihood**
1042        **Trees for Large Alignments**. *PLoS One* 2010, **5**(3):e9490.
1043 83.   Team RC: **R: A language and environment for statistical computing**. In: *R*
1044        *Foundation for Statistical Computing*. 3.3.2 edn; 2016.
1045 84.   Oksanen J, Blanchet FG, Kindt R, Legendre P, O'Hara RB, Simpson GL, Solymos P,
1046        Stevens MHH, Wagner H: **vegan: Community Ecology Package**. In., 1.17-6 edn:
1047        http://CRAN.R-project.org; 2011.
1048 85.   Legendre P, Anderson MJ: **Distance-based redundancy analysis: Testing**
1049        **multispecies responses in multifactorial ecological experiments**. *Ecological*
1050        *Monographs* 1999, **69**(1):1-24.
1051 86.   Anderson MJ: **A new method for non-parametric multivariate analysis of variance**.
1052        *Austral Ecology* 2001, **26**(1):32-46.
1053 87.   Hothorn T, Hornik K, Van de Wiel MA, Zeileis A: **A Lego system for conditional**
1054        **inference**. *American Statistician* 2006, **60**(3):257-263.
1055 88.   Kembel SW, Cowan PD, Helmus MR, Cornwell WK, Morlon H, Ackerly DD, Blomberg
1056        SP, Webb CO: **Picante: R tools for integrating phylogenies and ecology**.
1057        *Bioinformatics* 2010, **26**(11):1463-1464.
1058 89.   Paradis E, Claude J, Strimmer K: **APE: analyses of phylogenetics and evolution in R**
1059        **language**. *Bioinformatics* 2004, **20**(2):289-290.
1060 90.   Pinheiro J, Bates D, DebRoy S, Sarkar D, Team RDC: **nlme: Linear and Nonlinear**
1061        **Mixed Effects Models**. In.: http://CRAN.R-project.org; 2011.
1062 91.   Hommel G: **A stagewise rejective multiple test procedure based on a modified**
1063        **Bonferroni test**. *Biometrika* 1988, **75**(2):383-386.
1064 92.   Benjamini Y, Yekutieli D: **The control of the false discovery rate in multiple testing**
1065        **under dependency**. *Annals of Statistics* 2001, **29**(4):1165-1188.
1066 93.   Wang J, Thingholm LB, Skieceviciene J, Rausch P, Kummen M, Hov JR, Degenhardt F,
1067        Heinsen F-A, Ruhlemann MC, Szymczak S *et al*: **Genome-wide association analysis**
1068        **identifies variation in vitamin D receptor and other host factors influencing the gut**
1069        **microbiota**. *Nat Genet* 2016, **48**(11):1396-1406.
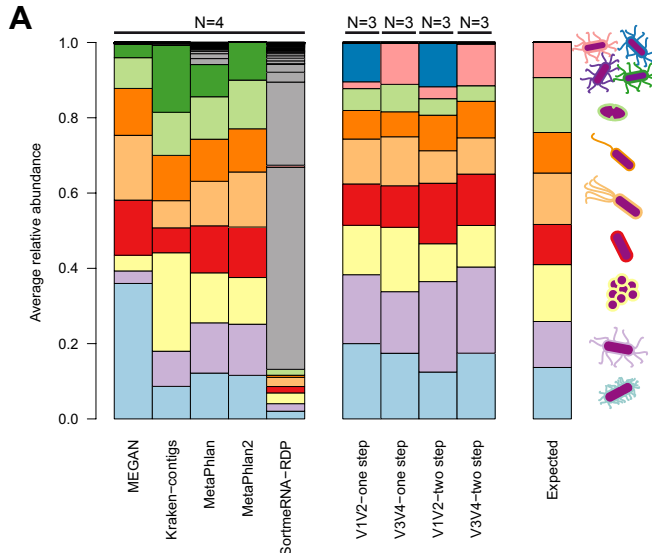
1070    94.    Rausch P, Rehman A, Künzel S, Häsler R, Ott SJ, Schreiber S, Rosenstiel P, Franke A,
1071           Baines JF: **Colonic mucosa-associated microbiota is influenced by an interaction**
1072           **of Crohn disease and FUT2 (Secretor) genotype**. *Proceedings of the National*
1073           *Academy of Sciences* 2011, **108**(47):19030-19035.
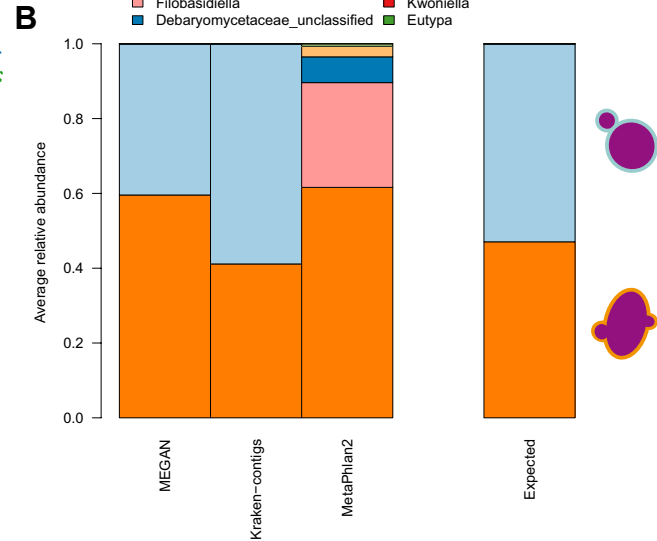
1074

**Bacteria:**
- Bacillus
- Salmonella
- Staphylococcus
- Lactobacillus
- Listeria
- Pseudomonas
- Enterococcus
- Escherichia
- Escherichia.Shigella
- Enterobacteriaceae
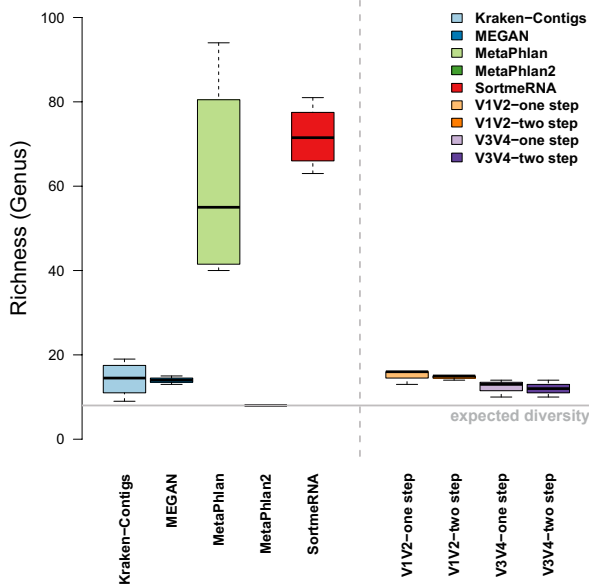- Shigella
- Others

**Fungi:**
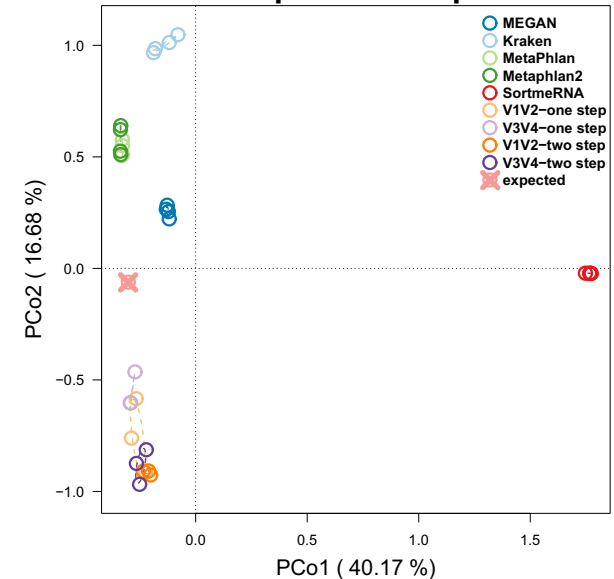- Saccharomyces
- Cryptococcus
- Filobasidiella
- Debaryomycetaceae_unclassified
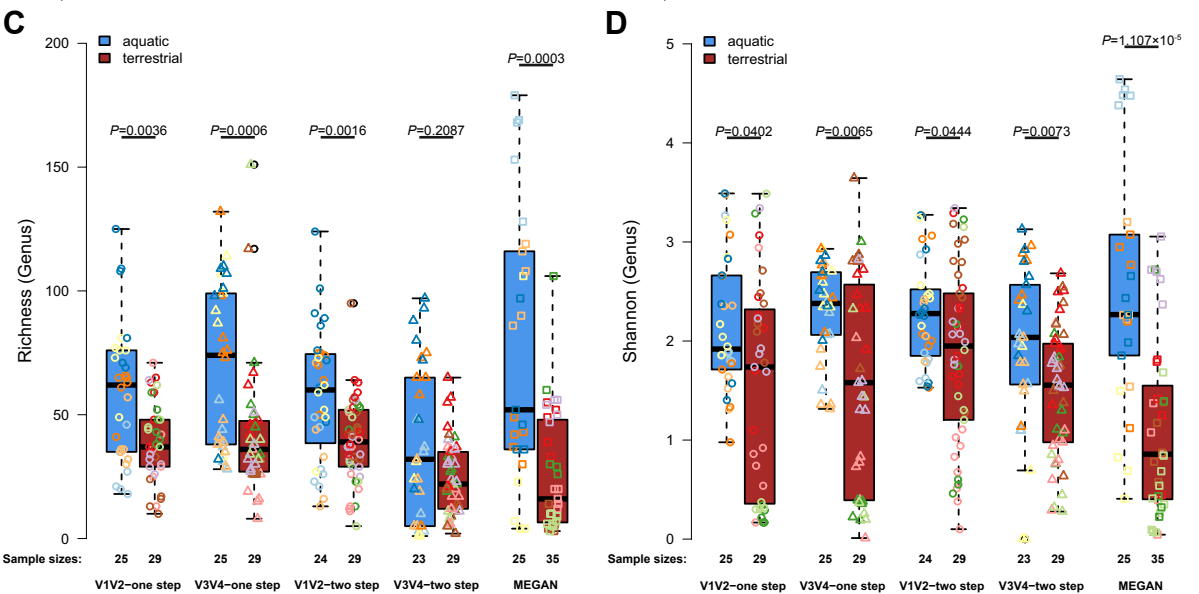- Naumovozyma
- Eremothecium
- Kwoniella
- Eutypa

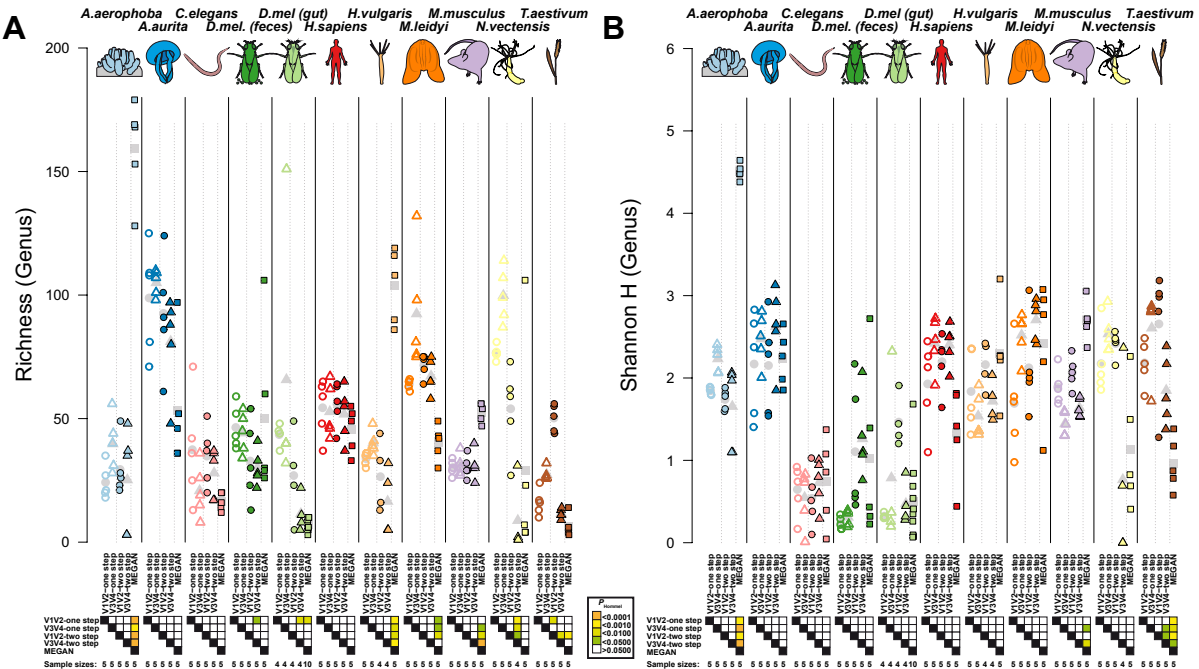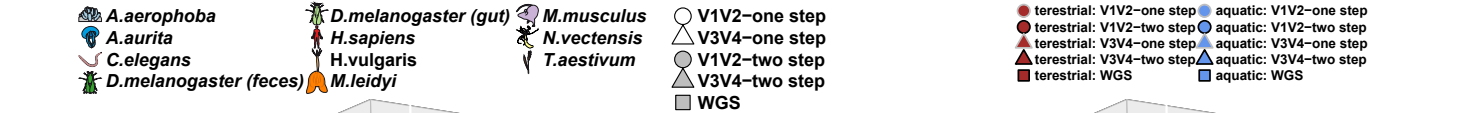**Genus distribution**

**A** — Average relative abundance. N=4 (MEGAN, Kraken−contigs, MetaPhlan, MetaPhlan2, SortmeRNA−RDP). N=3 (V1V2−one step, V3V4−one step, V1V2−two step, V3V4−two step). Expected.

**B** — Average relative abundance (MEGAN, Kraken−contigs, MetaPhlan2, Expected).

**C Community diversity**
Richness (Genus)
- Kraken−Contigs
- MEGAN
- MetaPhlan
- MetaPhlan2
- SortmeRNA
- V1V2−one step
- V1V2−two step
- V3V4−one step
- V3V4−two step

expected diversity

**D Distance to expected composition**
PCo2 ( 16.68 %)
PCo1 ( 40.17 %)
- MEGAN
- Kraken
- MetaPhlan
- Metaphlan2
- SortmeRNA
- V1V2−one step
- V3V4−one step
- V1V2−two step
- V3V4−two step
- expected

Legend:

*A.aerophoba*
*A.aurita*
*C.elegans*
*D.melanogaster (feces)*
*D.melanogaster (gut)*
*H.sapiens*
H.vulgaris
*M.leidyi*
*M.musculus*
*N.vectensis*
*T.aestivum*

V1V2−one step
V3V4−one step
V1V2−two step
V3V4−two step
WGS

terestrial: V1V2−one step     aquatic: V1V2−one step
terestrial: V1V2−two step     aquatic: V1V2−two step
terestrial: V3V4−one step     aquatic: V3V4−one step
terestrial: V3V4−two step     aquatic: V3V4−two step
terestrial: WGS               aquatic: WGS

Bray-Curtis distance

A
Stress: 0.17

B
Stress: 0.17

Jaccard distance

C
Stress: 0.144

D
Stress: 0.144

**Bray-Curtis**

**Jaccard**

Legend:
- ○ V1V2–one step
- ● V1V2–two step
- ■ MEGAN
- △ V3V4–one step
- ▲ V3V4–two step

- ● EggNOG-MEGAN
- ▲ EggNOG-Assembly
- ▲ CaZy-Assembly
- ○ EggNOG Cat.-MEGAN
- △ EggNOG Cat.-Assembly
- △ CaZy Cat.-Assembly

Y-axis: Procrustes correlation ($\sqrt{1-ss}$)

X-axis: P, C, O, F, G, S-OTU, Function WGS

**A** Gene richness (EggNOG-MEGAN) plotted for individual species: *A.aerophoba*, *A.aurita*, *C.elegans*, *D.melanogaster (feces)*, *D.melanogaster (gut)*, *H.sapiens*, *H.vulgaris*, *M.leidyi*, *M.musculus*, *N.vectensis*, *T.aestivum*.

Legend: *A.aerophoba*, *A.aurita*, *C.elegans*, *D.melanogaster (feces)*, *D.melanogaster (gut)*, *H.sapiens*, *H.vulgaris*, *M.leidyi*, *M.musculus*, *N.vectensis*, *T.aestivum*

**B** Gene richness (EggNOG-MEGAN) for aquatic vs terrestrial. *P*=0.6751

aquatic / terrestrial

**C** NMDS1, NMDS2, NMDS3. Stress: 0.074

**D** NMDS1, NMDS2, NMDS3. Stress: 0.074