

A primer and discussion on DNA-based microbiome data and related bioinformatics analyses

Gavin M. Douglas¹ and Morgan G. I. Langille^{2,*}

¹*Department of Microbiology and Immunology*

²*Department of Pharmacology, Dalhousie University, Halifax, Nova Scotia, Canada*

**Email for correspondence: morgan.langille@dal.ca*

Abstract

The past decade has seen an eruption of interest in profiling microbiomes through DNA sequencing. The resulting investigations have revealed myriad insights and attracted an influx of researchers to the research area. Many newcomers are in need of primers on the fundamentals of microbiome sequencing data types and the methods used to analyze them. Accordingly, here we aim to provide a detailed, but accessible, introduction to these topics. We first present the background on marker-gene and shotgun metagenomics sequencing and then discuss unique characteristics of microbiome data in general. We highlight several important caveats resulting from these characteristics that should be appreciated when analyzing these data. We then introduce the many-faceted concept of microbial functions and several controversies in this area. One controversy in particular is regarding whether metagenome prediction methods (i.e. based on marker gene sequences) are sufficiently accurate to ensure reliable biological inferences. We next highlight several underappreciated developments regarding the integration of taxonomic and functional data types. This is a highly pertinent topic because although these data types are inherently connected, they are often analyzed independently and primarily only linked anecdotally in the literature. We close by providing our perspective on this topic in addition to the issue of reproducibility in microbiome research, which are both crucial data analysis challenges facing microbiome researchers.

Background

Microbial communities encompass most of the genetic and species-level diversity on Earth. These communities are commonly characterized through DNA sequencing, which can be used to identify the presence and relative abundance of microbes in a community. These communities, including both the microbes, their constituent genes, and metabolites, are referred to as microbiomes. Due to technological improvements and the reduced cost of sequencing, the number of sequenced microbiomes has substantially grown in recent years. For instance, in 2017 the Earth Microbiome Project published a meta-analysis of 23,828 sequencing samples from all seven continents (Thompson et al. 2017). These data represented 109 environmental groupings and 21 major biomes, such as animal secretions, saline water, and soil. A key goal of microbial ecology research is to robustly analyze and correctly interpret these and other such microbial profiles.

But is DNA sequencing the best method for characterizing microbial communities? It is commonly observed that microbiome research would benefit from more emphasis on culturing, which enables individual microbes to be isolated and precisely studied in the lab. Traditionally, microbial communities were difficult to study by culturing alone because the vast majority of environmental microbes, particularly bacteria, could not be grown under standard culturing conditions (Staley and Konopka 1985). This issue remains unresolved even after gradual improvements to standard culturing conditions; a recent evaluation of six major environments identified only 34.9% of bacteria as culturable under standard conditions (Martiny 2019). However, modified culturing conditions can largely resolve this problem. By systematically applying 66 different conditions it was demonstrated that 95% of bacterial species in human stool samples could be grown in the lab (Lau et al. 2016). Therefore, it is no longer true for human stool samples, and likely other environments as well, that the majority of constituent bacteria cannot be cultured.

Despite these advances, a clear remaining advantage of DNA sequencing is that it enables microbial communities to be characterized in place, which theoretically enables the exact community relative abundances to be profiled. In practice, biases during sample collection and sequencing library preparation can perturb microbial relative abundances (Jones et al. 2015; Bukin et al. 2019; Watson et al. 2019). But nonetheless, DNA sequencing provides a more accurate view of the relative abundances of the com-

munity members than would be possible from culturing alone. For this reason, DNA sequencing remains the predominant method for characterizing microbial communities, although it is well-complemented by culturing (Lau et al. 2016).

DNA sequencing data is typically analyzed to identify specific associations between individual features (e.g. individual microbes) and sample groups of interest. Most commonly, researchers are interested in identifying associations between disease states and the relative abundance of features. A similar goal is often to investigate whether different measures of diversity in the studied dataset are associated with the sample groups. These measures of diversity are divided into alpha and beta diversity (Goodrich et al. 2014). Alpha diversity metrics refer to within-sample measures, such as richness (i.e. the number of taxa), and the Shannon diversity index (or entropy), which incorporates both the abundance and evenness of taxa within a sample (Jost 2006). In contrast, beta diversity refers to metrics that summarize variation between samples, which is most often performed by metrics that take the presence and abundance of features into account, such as the Bray-Curtis dissimilarity metric (Goodrich et al. 2014). Other microbiome-specific metrics have also been developed, such as the weighted UniFrac distance, which also takes the phylogenetic distance between taxa into account (Lozupone and Knight 2005). There is often more statistical power to detect overall differences based on alpha and beta diversity metrics than to detect associations with individual features, but diversity-level insights are also less actionable (Shade 2017).

There are many sub-categories of DNA sequencing approaches for characterizing microbial communities. One key distinction is between approaches that aim to characterize taxa (i.e. a group of organisms) and those that characterize genes and pathways, referred to as functions, that could be active in the community. These data types are referred to as taxonomic and functional microbiome data, respectively. Biologically this dichotomy is counter-intuitive; clearly genes are encoded in the genomes of taxa. So why does this distinction exist?

The reason is entirely related to methodological challenges. The most common and cost-effective sequencing approach focuses on sequencing marker genes. This method provides no direct information on the genomes of sequenced microbes, and instead is used to profile taxa. In contrast, shotgun metagenomics sequencing (MGS) provides information on all DNA present in a sample. MGS data can be used for analyzing both taxonomic and functional

132 profiles. However, it is difficult to integrate the
133 two data types, largely due to the complexity of
134 microbial communities and the fragmented nature
135 of DNA sequencing: it is relatively straight-forward
136 to identify genes in MGS data but challenging to
137 determine from which genomes they originated.

138 Herein we introduce the key forms of these data
139 types and highlight important caveats that should
140 be considered when they are analyzed. We first
141 cover the fundamentals of microbiome data analysis,
142 starting with marker-gene sequencing, and then move
143 to recently developed tools that could be leveraged to
144 conduct joint analyses of taxonomic and functional
145 data types. We conclude by highlighting two impor-
146 tant challenges that must be addressed in microbiome
147 data analysis.

148 **Marker-gene sequencing**

149 The earliest developed and most common form of mi-
150 crobiome sequencing is marker-gene sequencing, also
151 known as amplicon sequencing. Under this approach
152 specific genes are PCR-amplified and then sequenced.
153 There are two key requirements for robust marker
154 genes. First, they must be encoded by all (or at least
155 most) taxa of interest. Second, the observed sequence
156 divergence between orthologs should be approxi-
157 mately equal to the neutral mutation fixation rate
158 multiplied by double the divergence time between or-
159 thologs (Woese 1987). Note that the divergence time
160 should be doubled because mutations could accumu-
161 late in either lineage since the organisms diverged.
162 Genes displaying this second requirement have been
163 referred to as molecular chronometers. This term
164 highlights the close link between these marker genes
165 and the concept of the molecular clock (Zuckerlandl
166 and Pauling 1965): given equal mutation rates and
167 equal fixation rates for neutral mutations, the number
168 of neutral substitutions between organisms is directly
169 proportional to the evolutionary divergence between
170 them.

171 However, there are many reasons why a gene
172 might be an unreliable molecular chronometer (Janda
173 and Abbott 2007). One reason is that if a gene varies
174 in function across taxa then contrasting selection
175 pressures could result in different non-synonymous
176 substitution rates (Wheeler et al. 2016). For
177 instance, as previously observed (Woese 1987), the
178 cytochrome complex gene is a useful molecular
179 chronometer in eukaryotes, but suffers from draw-
180 backs. This gene was shown to be useful for building
181 early phylogenetic trees that represented both long
182 evolutionary distances across eukaryotes and short

183 distances between human populations (Fitch and
184 Margoliash 1967). However, within prokaryotes the
185 cytochrome complex systematically varies in size,
186 which is believed to be due to positive selection (Am-
187 bler et al. 1979). Because positive selection is likely
188 driving divergence between orthologous cytochrome
189 complexes, in at least some cases it would be an in-
190 valid molecular chronometer to study in prokaryotes.
191 Similarly, if a gene is sufficiently divergent between
192 organisms then it can be difficult to accurately align
193 residues. Misalignments lead to inaccurate estimates
194 of evolutionary divergence, which is particularly true
195 if the gene accumulates insertions and deletions. Such
196 highly divergent regions, particularly in areas under
197 no selective constraint, have been referred to as
198 “evolutionary stopwatches” (Woese 1987), because
199 they are useful only at short evolutionary distances.
200 Therefore, to select a robust marker gene one should
201 adhere in some ways to the Goldilocks principle: some
202 nucleotide conservation is needed, but not too much.

203 **The 16 Svedberg (16S) ribosomal RNA (rRNA)**
204 **gene fits well with this principle.** This gene features
205 highly conserved regions surrounding nine less con-
206 served regions (referred to as variable regions). It
207 is also encoded by all prokaryotes and represents 50
208 helical RNA regions encoded by approximately 1,500
209 base-pairs (Woese et al. 1980). This high number
210 of independent functional domains is valuable in a
211 marker gene (Woese 1987). This is because if there
212 are non-random substitutions within a single domain,
213 but substitutions in the majority of other domains
214 are driven by random processes, there would likely be
215 little effect on estimates of evolutionary divergence.
216 This gene also encodes a highly conserved function
217 across both prokaryotes and eukaryotes (where it is
218 called the 18S rRNA gene). The 16S rRNA molecule
219 is part of the 30S small subunit (SSU) of the ribo-
220 some, which helps initiate protein synthesis by bind-
221 ing the Shine-Dalgarno sequence in messenger RNA
222 (mRNA) to align the ribosome with the encoded start
223 codon. Many changes in the highly conserved region
224 of the 16S rRNA gene affect its binding affinity to the
225 ribosome and mRNA. The strong negative selection
226 acting against such substitutions makes these regions
227 valuable for detecting rare substitutions between
228 distant relatives, anchoring alignments of 16S rRNA
229 genes, and for primer design (Wang et al. 2013).

230 Since the 16S rRNA gene was identified as a useful
231 molecular chronometer, it has been the prime marker
232 gene used to develop phylogenetic models of the tree
233 of life. Most famously, an alignment of 16S (and
234 18S) rRNA gene sequences from across life lead to
235 distinguishing archaea, bacteria, and eukaryotes into
236 distinct domains (Woese and Fox 1977). In these

237 early days, research focused on analyzing the rRNA
238 sequences of isolated microbes. This was painstaking
239 work, as illustrated by the prediction in 1987 that
240 future research groups could plausibly sequence on
241 the order of one hundred 16S rRNAs a year (Woese
242 1987).

243 Thirty-four years later, through next-generation
244 sequencing technology, insufficient availability of se-
245 quenced rRNA genes is no longer a common com-
246 plaint. Databases such as SILVA contain enormous
247 collections of sequenced SSU fragments; as of Au-
248 gust 2020 SILVA contained 9,469,124 non-clustered,
249 independent sequences (Quast et al. 2013). Software
250 such as redbiom also enables unique 16S rRNA gene
251 variants to be compiled from the growing number
252 of 16S rRNA gene sequencing (hereafter referred to
253 as 16S sequencing) studies (McDonald et al. 2019).
254 These 16S datasets are produced to characterize
255 and compare the relative abundances of prokaryotes
256 across communities. However, despite the ubiquity
257 of such datasets, they are non-trivial to analyze
258 and interpret. There are numerous methodological
259 reasons for this difficulty.

260 First, due to sequencing length constraints, only
261 certain 16S rRNA gene variable regions are typically
262 amplified and sequenced. Each variable region has
263 particular strengths and limitations (Chen et al.
264 2019; Johnson et al. 2019). Along with our colleagues
265 we have previously compared the biases between the
266 amplified fragments from variable regions four and
267 five and from regions six to eight (written as V4-V5
268 and V6-V8, respectively) on a mock community from
269 the Human Microbiome Project (HMP) (Comeau et
270 al. 2017). We found the V4-V5 region overrep-
271 resented Firmicutes and Bacteroides while drastically
272 underestimating Actinobacteria. In contrast, the
273 V6-V8 region overrepresented Proteobacteria and
274 underrepresented Bacteroides. These biases highlight
275 that choice of variable region can depend on which
276 taxa are of interest. For example, region V4-V5
277 was recently shown to be superior to region V6-V8
278 for identifying archaea in the North Atlantic Ocean
279 (Willis et al. 2019). In this case the authors were
280 particularly interested in archaeal diversity so the V4-
281 V5 region was more appropriate.

282 Typically, however, the taxonomic scope of inter-
283 est and region biases in a particular environment are
284 not clear and little or no rationale is given for the
285 variable region selection. This is a problem, because
286 analyses of the same communities with different
287 variable regions can result in not only systematic
288 biases in the raw data, but also in strikingly different
289 biological interpretations. For example, key species
290 that modulate human vaginal health are underrep-

resented or missing in V1-V2 sequencing datasets, 291
such as *Gardnerella vaginalis*, *Bifidobacterium bi-* 292
fidum, and *Chlamydia trachomatis* (Graspeuntner et 293
al. 2018). Application of this region for profiling 294
vaginal samples, instead of the more appropriate 295
choice of the V3-V4 region, can result in entirely 296
missing associations between vaginal health and the 297
microbiome. Similarly, a comparison of the tick 298
microbiome based on six sequenced 16S rRNA gene 299
regions found a wide range of the number of prokary- 300
otic families and in the Shannon diversity index for 301
each individual tick (Sperling et al. 2017). The 302
problem of such biases in variable region selection 303
is beginning to recede as long-read technologies, such 304
as that developed by Pacific Biosciences of California 305
and Oxford Nanopore Technologies Limited, enable 306
full-length 16S sequencing (Callahan et al. 2019; 307
Johnson et al. 2019). However, it will remain an 308
important issue for the foreseeable future as long 309
as the microbiome is largely studied by short-read 310
sequencing. 311

Regardless of the sequenced region, most reads 312
originating from the same biological molecule will 313
differ due to sequencing errors. Raw reads are either 314
clustered based on sequence identity into operational 315
taxonomic units (OTUs) or alternatively errors are 316
corrected to produce amplicon sequence variants 317
(ASVs). OTUs are typically clustered at 97% identity 318
(Goodrich et al. 2014), which often results in merging 319
different species into a single OTU (Mysara et al. 320
2017). This issue has long plagued 16S rRNA gene- 321
based analyses. For instance, *Bacillus globisporus* 322
and *Bacillus psychrophilus* are problematic cases be- 323
cause their 16S genes share 99.5% sequence identity, 324
but are highly distinct at the genome level (Fox et al. 325
1992). 326

327 In contrast to clustering approaches, error-
328 correcting approaches, referred to as denoising meth-
329 ods, theoretically can correct raw reads sufficiently
330 well to produce exact biological molecules. Several
331 different denoising approaches have emerged recently.
332 DADA2 is the most sophisticated approach, which
333 generates a different parametric error model for every
334 input sequencing dataset (Callahan et al. 2016a).
335 The raw sequencing reads are then corrected to
336 generate ASVs based on this error model. Deblur
337 and UNOISE3 are two other denoising tools that are
338 based on rapidly clustering raw reads and using pre-
339 determined hard cut-offs related to the expected error
340 rates to generate ASVs. We and other colleagues
341 have evaluated the performance of these three tools
342 and open-reference OTU clustering (which combines
343 both de novo and reference-based clustering) and
344 found that all three denoising methods result in

345 similar overall microbial communities (Nearing et al.
346 2018). In contrast, we found that open-reference
347 OTU clustering resulted in a high rate of spurious
348 OTUs compared to these methods. Nonetheless,
349 there were important differences between the three
350 methods, particularly in terms of richness and when
351 profiling rare taxa (Nearing et al. 2018). A more
352 recent independent validation based on a higher
353 number of test datasets reached similar conclusions
354 (Prodan et al. 2020).

355 In addition to 16S rRNA gene sequencing data,
356 there are multiple marker genes appropriate for profil-
357 ing eukaryotic diversity. As mentioned above, the 18S
358 rRNA gene is the homolog of the 16S rRNA gene in
359 eukaryotes and is widely used to profile that domain.
360 However, fungi are more difficult to distinguish based
361 on the 18S rRNA gene, because fungi lack several
362 variable regions for this gene (Schoch et al. 2012).
363 Instead, the internal transcribed spacer (ITS) region,
364 although not strictly a marker gene, is more often
365 amplified to study fungal communities, because it
366 typically has more resolution to distinguish fungi
367 than the 18S rRNA gene (Liu et al. 2015). This
368 region is within the nuclear rRNA cistron of fungi
369 genomes, which contains the 18S, 5.8S, and the 28S
370 rRNA genes. The ITS regions encompasses the two
371 intergenic regions, which have relatively high rates
372 of insertions and deletions, and the 5.8S rRNA gene
373 (Schoch et al. 2012). Only a single intergenic
374 region is typically amplified, referred to as ITS1 or
375 ITS2, which have better discriminatory resolution
376 for the major phyla Basidiomycota and Ascomycota,
377 respectively (Saroj et al. 2015).

378 Although the marker genes described above are
379 the most commonly profiled loci, in many cases
380 there are marker genes more appropriate for specific
381 lineages. For example, several halophilic species of
382 *Haloarcula* encode multiple 16S copies that can differ
383 by more than 5% sequence identity within the same
384 genome (Sun et al. 2013). Consequently, different
385 marker genes are often used when building phylo-
386 genetic trees representing a single species or genera.
387 The chaperonin-60 (*cpn60*) gene is one useful alter-
388 native prokaryotic marker gene, which is particularly
389 useful for distinguishing taxa at resolutions below the
390 genus level (Links et al. 2012). For example, the
391 *cpn60* gene has been frequently profiled in vaginal
392 microbiome samples, because variation at this locus
393 can distinguish subgroups of *Gardnerella vaginalis*
394 that cannot be distinguished based on the 16S rRNA
395 gene alone (Jayaprakash et al. 2012). More generally,
396 marker genes for specialized comparisons are often
397 chosen to match the defining function of a given
398 lineage. For example, the methyl coenzyme M re-

dundance A (*mrcA*) gene and a nitrate reductase gene
have been previously profiled to explore the diversity
of methanogens (Hallam et al. 2003) and nitrogen-
fixing microbes (Comeau et al. 2019), respectively.

Shotgun metagenomics sequencing

Shotgun metagenomics sequencing (MGS) is a quali-
tatively different method from marker-gene sequenc-
ing, because it involves sequencing all DNA in a
community. This is a major advantage and means
that MGS data can profile any taxa, including viruses
and microbial eukaryotes. MGS approaches were first
applied to study ocean water communities through
a Fosmid cloning approach (Stein et al. 1996).
Building upon such early studies, the potential for
leveraging MGS was widely publicized by an investi-
gation into the microbial diversity of the Sargasso
Sea (Venter et al. 2004). This study identified
1.2 million previously unknown genes and many
other microbial features that would be impossible
to study with 16S rRNA gene sequencing. These
and other related observations sparked an explosion
of interest in profiling microbial communities with
MGS approaches. This interest has culminated in
the generation of enormous MGS datasets such as
the ongoing work on the Earth Microbiome Project
(Thompson et al. 2017) and the Human Microbiome
Project (Lloyd-Price et al. 2017).

There are two main approaches for analyzing
MGS data: read-based workflows and metagenomics
assembly. Each of these approaches has strengths
and weaknesses, but in both cases the generated
profiles imprecisely reflect biological reality. For
instance, the number of species identified by different
read-based methods can vary by three orders of
magnitude (McIntyre et al. 2017). The exact species
relative abundances can also drastically differ across
tools, as recently shown in a comparison of read-
based methods applied to simulated datasets (Ye et
al. 2019). Different approaches for metagenomic
assembly will produce different assembled contigs
and microbial profiles as well (Olson et al. 2019).
Unsurprisingly, given this wide variation, there is also
low concordance between 16S sequencing and MGS
data taken from the same samples. For example,
one comparison found that fewer than 50% of phyla
identified in water samples based on 16S sequencing
were also identified in the corresponding MGS profiles
(Tessler et al. 2017). This wide variation in results
highlights that any interpretation of MGS profiles,
similar to 16S profiles, should be done cautiously. It

450 is crucial to appreciate that any approach will have
451 important weaknesses and that the generated profile
452 will only partially represent the actual microbial
453 diversity.

454 With those important caveats in mind, an under-
455 standing of the different approaches is nonetheless
456 important to give context to MGS data analysis.
457 Read-based workflows involve little or no assem-
458 bly of the reads and instead each read (or pair
459 of reads) is treated independently. This is the
460 most common approach for analyzing MGS data,
461 particularly because it can be performed with low
462 sequencing depth (Hillmann et al. 2018) and in
463 complex communities (Zhou et al. 2015). However,
464 an important disadvantage of this approach is that
465 taxonomic and functional annotations are typically
466 generated and treated as entirely independent data
467 types (Figure 1a). It is also possible to map reads
468 against a set of known reference genomes, which does
469 link the two data types (Figure 1b). Although this
470 is an invaluable approach when applied to genomes
471 assembled from the study environment (see below),
472 the results are typically near incomprehensible when
473 reads are mapped against a database of thousands
474 of genomes. Instead, the most common approach
475 for generating taxonomic profiles is either based on a
476 marker-gene or k-mer method.

477 Marker-gene approaches are based on the insight
478 that specific genes can be used to identify the pres-
479 ence and relative abundance of certain taxa. An
480 extreme example is to use solely the 16S rRNA
481 gene for taxonomic classification (Hao and Chen
482 2012). More commonly, marker-gene approaches
483 base classifications on many genes. For instance,
484 PhyloSift (Darling et al. 2014) leverages 37 nearly
485 universal prokaryotic marker-genes (Wu et al. 2013)
486 in addition to eukaryotic and viral gene sets to make
487 a combined set of approximately 800 (mainly viral)
488 gene families for classification. Aligned reads are
489 placed into a phylogenetic tree of reference sequences
490 and taxonomic classification is performed based on
491 summing the likelihood of each taxa based on each
492 read placement (Darling et al. 2014). MetaPhlAn2 is
493 a contrasting approach that instead bases taxonomic
494 predictions on the presence of clade-specific marker
495 genes, which are genes only found in that given
496 lineage, and found in all members (Truong et al.
497 2015). This method has rapidly become the most
498 popular marker-gene MGS approach. However, given
499 that this approach is limited by the existence of
500 robust clade-specific genes, it is not surprising that
501 it tends to have low sensitivity (Tessler et al. 2017;
502 Miossec et al. 2020), meaning that it misses taxa that
503 are actually present.

504 In contrast, k-mer-based approaches are much
505 more sensitive but have slightly lower specificity
506 than marker-gene methods (Miossec et al. 2020).
507 These approaches search for exact matches of short
508 DNA sequences (k-mers) within reference genomes.
509 An algorithm such as lowest-common ancestor is
510 then performed to determine the likely taxonomic
511 classification based on all matching genomes. Two
512 common kmer-based approaches are kraken2 (Wood
513 et al. 2019) and centrifuge (Kim et al. 2016), both of
514 which match k-mers against a compressed database
515 of reference genomes. In contrast to the marker-gene
516 results, the main challenge of analyzing taxonomic
517 profiles output by these methods is the high number
518 of rare taxa of different ranks identified, some of
519 which may be false positives. Summarizing the
520 output profiles with an additional approach, such as
521 the Bayesian abundance re-estimation tool Bracken
522 (Lu et al. 2017) in the case kraken2 data, can help
523 partially mitigate this problem.

524 Most functional read-based methods are based
525 on a similarity search of reads against a database
526 of known gene families. This is primarily done in
527 protein space, because protein similarity matches are
528 more informative and the database requirements are
529 lower (Koonin and Galperin 2003). The common sim-
530 ilarity searching tool BLASTX is prohibitively slow
531 when scanning millions of reads, which has driven the
532 development of faster alternatives like DIAMOND
533 (Buchfink et al. 2015) and MMseqs2 (Steinegger
534 and Söding 2017). These faster alternatives are
535 leveraged by workflows implemented in software such
536 as MEGAN (Huson et al. 2007) and HUMAnN2
537 (Franzosa et al. 2018) to identify gene family matches
538 and output overall metagenome profiles. HUMAnN2
539 is a unique approach in that it first screens reads that
540 map to reference genomes of taxa identified as present
541 with MetaPhlAn2. This step enables a small subset of
542 gene families to be linked directly to particular taxa.
543 However, the vast majority of gene families typically
544 have no taxonomic links and are only part of the
545 community-wide metagenome. There are clear issues
546 with the general approach implemented by these gene
547 profiling approaches, as has been previously observed:
548 “genes are expressed in cells, not in a homogenized
549 cytoplasmic soup” (McMahon 2015).

550 Linking functional annotations to specific taxa
551 by assembling raw reads is the ideal approach to
552 resolve this problem, but this too comes with caveats.
553 Most importantly, insufficiently high read depth,
554 which depends on the complexity of a sample, can
555 result in too few assembled contigs to sensibly ana-
556 lyze. Nonetheless, with sufficiently high read depth
557 metagenome assembly can be a valuable way to

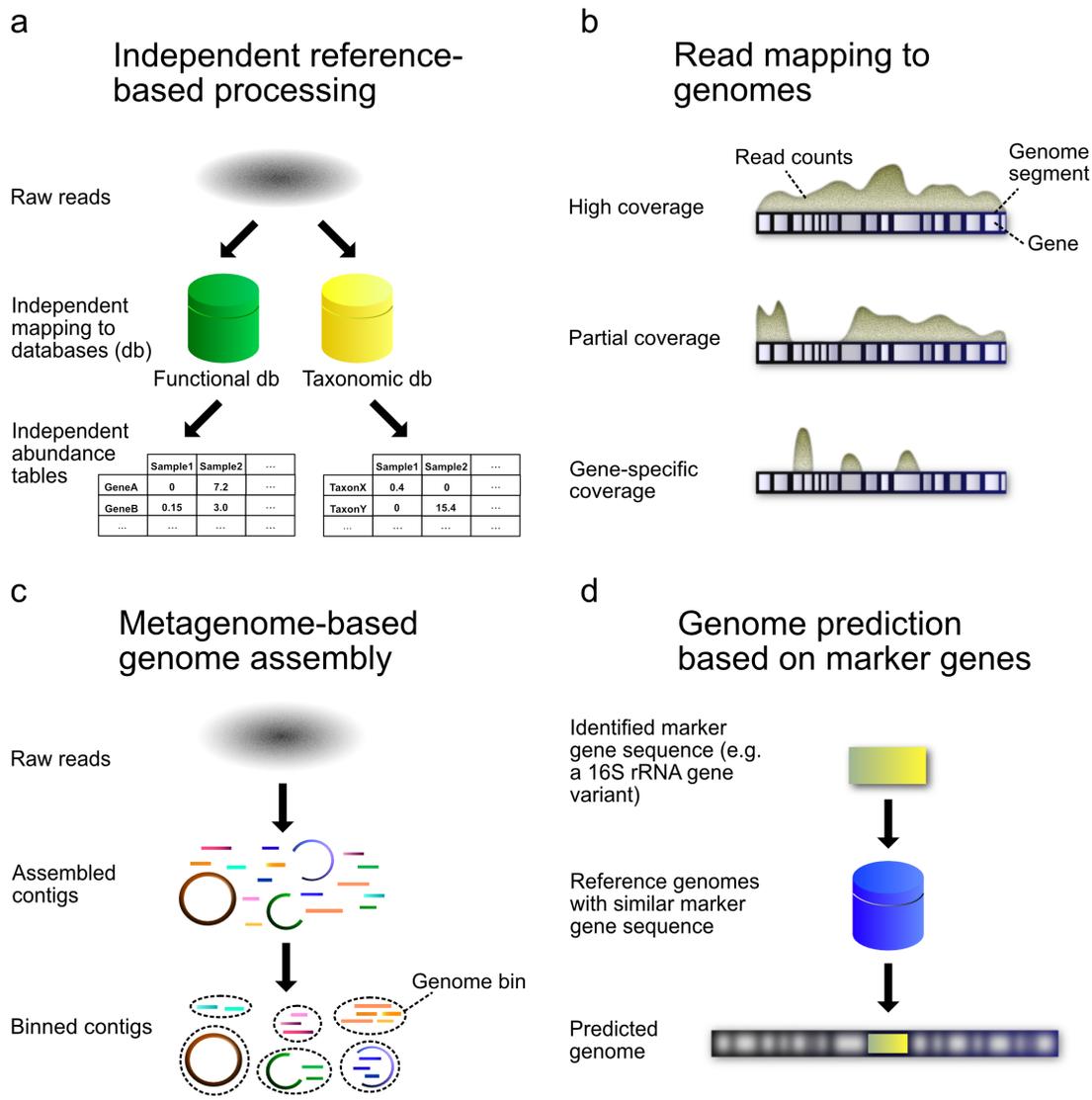


Figure 1: Key approaches for generating joint taxonomic and functional data from microbiome sequencing data. (a) Read-based processing of shotgun metagenomics data to generate functional and taxonomic abundance tables independently. (b) Read mapping to genome sequences can be used to infer the presence of a taxon based on read coverage. It can also be used to identify the presence of strains missing specific genes or of the inverse: a community containing specific genes from a genome while the rest of the genome is absent. Note that all of these inferences are best made in low complexity communities where there are few ambiguous read mappings, and where the possible set of genomes present is relatively well defined. This is particularly applicable when mapping reads against metagenome assembled genomes from the same dataset. (c) Metagenomics-based genome assembly involves assembling reads into contigs and then binning contigs into categories representing metagenome-assembled genomes. Missing from this diagram is the important quality control step, which is essential to follow-up metagenomics assembly. Also, this approach is best for profiling dominant organisms, and produces the best results when sequencing read depth is high and/or community complexity is low. (d) Genome prediction based on marker gene sequences is another method of producing joint taxa and function profiles, which in this case are explicitly linked, similar to assembling genomes. However, these approaches are highly biased towards the specific reference genomes used for prediction. In addition, they can only predict genome content to the level at which the chosen marker gene differs between closely related taxa. This is a major limitation as many strains of bacteria with highly divergent genome content have identical marker gene sequences.

558 leverage information about microbial communities
559 (Figure 1c). There are many metagenome assembly
560 tools available, such as MetaSPAdes (Nurk et al.
561 2017) and Megahit (Li et al. 2015). The resulting
562 assembled contigs from these approaches are typically
563 categorized (or “binned”) into groups of contigs with
564 similar characteristics. This binning is primarily
565 performed by identifying contigs that are found at
566 similar relative abundances across samples and/or
567 that contain similar proportions of different k-mers
568 (particularly 3-mers) (Ayling et al. 2019). These bins
569 represent metagenome-assembled genomes (MAGs)
570 that must undergo stringent checks to help evaluate
571 the overall quality (Bowers et al. 2017). The
572 key method for performing quality control on these
573 genomes is to scan for known universal single-copy
574 genes (USCGs), with a tool such as CheckM (Parks
575 et al. 2015). The percentage of USCGs present
576 provides an estimate of overall genome completeness.
577 In contrast, the number of USCGs found in multiple
578 copies can be used to calculate the redundancy,
579 which is potential evidence for contamination in the
580 genome.

581 Characteristics of microbiome 582 count data

583 Regardless of the sequencing technology and work-
584 flow used for profiling a microbial community, the
585 final product is typically a count table. This is true
586 for many sequencing approaches, such as RNA se-
587 quencing, but there are several important differences.
588 First, unlike in the case of RNA sequencing where
589 there are typically a known number of genomic loci,
590 novel taxa and functions are frequently identified in
591 microbiome data. For instance, novel OTUs, ASVs,
592 and contigs are frequently identified in taxonomic
593 analyses. Similarly, 25-85% of proteins in MGS are
594 novel microbial genes of unknown function (Prakash
595 and Taylor 2012). Second, no statistical distribution
596 fits microbiome data in all contexts. For example,
597 many statistical distributions, including the negative
598 binomial (Love et al. 2014), beta binomial (Martin
599 et al. 2020), and Poisson (Faust et al. 2012)
600 distributions have been proposed as appropriate fits
601 to microbiome data. However, upon analysis with
602 real data these and other distributions fit with incon-
603 sistent accuracy (Weiss et al. 2017; Calgaro et al.
604 2020). Last, microbiome count tables typically have
605 high sparsity, meaning that there is a high proportion
606 of features not found across many samples (Thorsen
607 et al. 2016). These characteristics make microbiome
608 data analysis challenging for all taxonomic analyses

and most functional analyses (see Microbial functions
section).

609
610
611 These challenges are exacerbated by the inherent
612 compositionality of sequencing data. Compositional
613 data refers to data that is constrained to an arbi-
614 trary constant sum (Aitchison 1982), such as the
615 arbitrary number of raw sequencing reads output per
616 sample. This characteristic means that the observed
617 abundance of any given feature is dependent on the
618 observed abundance of all other features. A simple
619 example can help illustrate the implications of this
620 characteristic. Imagine a microbe, microbe x , at low
621 relative abundance in sample a and at high relative
622 abundance in sample b . An observer might naively
623 infer that there is more of microbe x in sample b than
624 in sample a . However, there are many reasons this
625 could be false. For instance, the absolute abundance
626 of microbe x could be the same in each sample but
627 the abundance of other microbes in general might be
628 higher in sample a . This higher total microbial load
629 would push the relative abundance of microbe x in
630 sample a down. Depending on the total microbial cell
631 count it is even possible that the absolute abundance
632 of microbe x could be higher in sample a than
633 in sample b , but that the relative abundance is
634 simply lower. This example highlights a necessary
635 consideration regarding microbiome sequencing data
636 analysis: it only provides information on the relative
637 abundances, or percentages, of features and does not
638 provide insight on feature absolute abundances.

639 This important characteristic was not widely ap-
640 preciated in the field until relatively recently, when
641 researchers identified fatal issues with common ap-
642 proaches for analyzing microbiome data (Gloor et
643 al. 2016, 2017). Standard differential abundance ap-
644 proaches, such as the t-test and Wilcoxon test, when
645 applied to relative abundances, and microbiome-
646 specific tools such as LEfSe (Segata et al. 2011)
647 do not account for this compositionality. Com-
648 mon summary metrics for microbiome data, such as
649 the UniFrac distance, also suffer from this problem
650 (Gloor et al. 2017). This is a major issue, because
651 ignoring this characteristic is known to lead to spu-
652 rious discoveries with compositional data (Aitchison
653 1982; Jackson 1997; Fernandes et al. 2014).

654 Fortunately, there is active work in the field
655 to resolve this issue and numerous compositional
656 approaches have been developed. The focus has
657 primarily been on developing novel correlation (Fried-
658 man and Alm 2012; Kurtz et al. 2015; Schwager
659 et al. 2017) and differential abundance approaches,
660 such as ALDEx2 (Fernandes et al. 2013, 2014) and
661 ANCOM (Mandal et al. 2015). A common theme
662 of these compositional approaches is that the data

663 is transformed based on the ratio of feature relative
664 abundances to some reference frame (Aitchison 1982;
665 Morton et al. 2019). This choice of reference
666 frame varies substantially between approaches. For
667 instance, ALDEx2 transforms relative abundances by
668 the centred log-ratio (CLR) transformation (Fernan-
669 des et al. 2013), which essentially normalizes feature
670 relative abundances by the mean relative abundance
671 per sample. This approach transforms the original
672 data but maintains the interpretation of individual
673 features. In contrast, it has been suggested that
674 analyses could instead be based on ratios between
675 features (Morton et al. 2019), which converts the
676 data type into comparisons of features rather than
677 individual features.

678 There are no best-practices regarding approaches
679 that compositionally transform individual features.
680 More generally, differential abundance tests com-
681 monly produce widely different sets of significant taxa
682 from each other (Thorsen et al. 2016; Weiss et al.
683 2017; Hawinkel et al. 2019). This wide variation is
684 largely due to specific characteristics of microbiome
685 count data. A large proportion of the variation in re-
686 sults is driven by high false discovery rates. Although
687 many methods advertise that only approximately 5%
688 of significant taxa are likely false positives, it has
689 been estimated that for some methods the actual false
690 discovery rate is substantially higher (Hawinkel et al.
691 2019). This particular validation observed this trend
692 for several methods, including ANCOM (Mandal et
693 al. 2015) and metagenomeSeq (Paulson et al. 2013),
694 two microbiome-oriented methods that are otherwise
695 considered conservative (Paulson et al. 2013; Weiss
696 et al. 2017). In addition, a recent evaluation of
697 differential abundance tools found that compositional
698 methods are actually less robust than several non-
699 compositional alternatives (Calgaro et al. 2020).

700 Given this wide variation in differential abun-
701 dance tool performance and unclear best-practices,
702 how is a microbiome researcher to proceed? One
703 possible answer is that a change in expectations
704 regarding the interpretability of microbiome data
705 analysis is needed. In particular, analyses using
706 ratios between the relative abundances of taxa has
707 been shown to be robust, although it comes at
708 the cost of interpretability (Morton et al. 2019).
709 However, an important issue is how to determine
710 which taxa should be the numerator and denominator
711 of each ratio. One solution is to leverage the
712 bifurcating structure of a clustered tree (Egozcue
713 and Pawlowsky-Glahn 2011; Morton et al. 2017)
714 or phylogenetic tree (Silverman et al. 2017) of
715 features. Analyses can be focused on the ratios in
716 relative abundances between features on the left-hand

717 and right-hand of each node in the tree. Despite
718 the potential of this approach, it is rarely used for
719 standard microbiome analyses because it is unclear
720 how to biologically interpret any differences in the
721 values of these ratios across samples.

722 This discussion of microbiome data characteristics
723 has focused on taxonomic features based on either
724 16S sequencing or read-based MGS data analysis.
725 However, it is important to emphasize that count
726 tables produced from MAGs do not resolve this issue.
727 In fact, attempting to account for these challenging
728 characteristics of microbiome count data and the
729 links between taxa and function makes the analysis
730 more difficult.

731 **Microbial functions**

732 To this point we have only discussed functional micro-
733 biome data in vague terms as referring to microbial
734 gene abundances. When based on DNA sequencing
735 data this information summarizes the functional po-
736 tential, meaning the functions that are present, but
737 not necessarily active in a community. However,
738 rather than individual gene sequences, research is
739 typically focused on gene families, which are gene
740 clusters. Alternatively, the focus is sometimes on
741 higher-order functional categories like pathways. To
742 complicate matters further, there are several different
743 functional ontologies, which are different frameworks
744 for studying functions at different resolutions. De-
745 pending on which of these functional ontologies and
746 sub-categories are analyzed, the characteristics of the
747 data can drastically differ.

748 The Universal Protein Resource (UniProt) Refer-
749 ence Clusters (UniRef) database contains all protein
750 sequences from the Swiss-Prot (manually curated)
751 and TrEMBL (automated) databases clustered at
752 either 50%, 90%, or 100% identity (Apweiler et al.
753 2004). The most recent versions of these clusters
754 have been generated with the MMseqs2 algorithm
755 (Steinegger and Söding 2018). As of June 30th, 2020,
756 the 100% identity clusters (called UniRef100), cor-
757 responded to 235,561,514 unique protein sequences,
758 which provides a detailed summary of almost all
759 known protein sequences. Despite being clustered
760 at lower identity thresholds, UniRef50 and UniRef90
761 nonetheless contain enormous numbers of protein
762 clusters: 41,883,832 and 115,885,342, respectively.

763 The UniRef database contrasts with another com-
764 mon functional ontology, the Kyoto Encyclopedia of
765 Genes and Genomes (KEGG) database (Kanehisa
766 and Goto 2000; Kanehisa et al. 2016). KEGG
767 is based on 23,530 individual gene families (as of

768 September 10th, 2020), which are called KEGG
769 orthologs (KOs). The advantage of KOs is that
770 the majority have well-described molecular functions
771 that can be linked to higher-order KEGG pathways
772 and modules. Accordingly, any analysis of KEGG
773 data will likely result in less sparse count tables than
774 the corresponding UniRef-based database, simply be-
775 cause KOs are shared across more taxa than UniRef
776 clusters.

777 To illustrate this point, we and our colleagues
778 have previously compared the taxonomic coverage of
779 each function within these two functional ontologies
780 and each sub-category (Inkpen et al. 2017). We
781 found that all UniRef functions, including those in
782 UniRef50 clusters, are on average found in a single
783 domain and encoded by fewer than four species. In
784 contrast, we found that KOs were encoded in 1.3
785 domains and 184.3 species on average. Similarly,
786 the high-level KEGG modules and pathways were
787 predicted to be potentially active in a mean of 1.7 and
788 2.5 domains and 671 and 1267.6 species, respectively
789 (Inkpen et al. 2017). Based on these statistics,
790 clearly a shift in the abundance of a UniRef cluster
791 should not be treated the same as a KEGG function:
792 the former corresponds to the activity of a small
793 number of species while the latter could correspond
794 to a large assemblage. This example highlights that
795 the choice of how function is defined in a given
796 analysis can have profound effects on the biological
797 interpretation.

798 In addition to UniRef and KEGG, several other
799 functional ontologies have been leveraged for micro-
800 biome analyses. Key examples of additional func-
801 tion types include: Clusters of Orthologous Genes
802 (COGs) (Tatusov et al. 2000; Makarova et al. 2015),
803 Enzyme Commission (EC) numbers, Protein families
804 (Pfam) (Punta et al. 2012; Finn et al. 2014), and
805 TIGRFAMs (Haft et al. 2003). These categories
806 represent a range of approaches for defining gene
807 families and functional categories.

808 The COG strategy for functional annotation
809 was originally intended to phylogenetically classify
810 proteins into groups of orthologs (Tatusov et al.
811 2000). This one-to-one approach of matching indi-
812 vidual orthologs has now been expanded to allow
813 for more complex relationships between genes, such
814 as paralogs and horizontally transferred homologs
815 (Makarova et al. 2015; Galperin et al. 2019). As of
816 2015, there were 4,631 independent COGs (Galperin
817 et al. 2015). The COG framework is similar to that
818 of the eggNOG database (Jensen et al. 2008), which
819 is a more high-throughput, automated approach.
820 However, the key advantage of the COG database
821 is that orthologous genes are clustered into 26 inter-

822 pretable functional categories, which are expanded
823 from categories originally defined to functionally bin
824 *Escherichia coli* genes (Riley 1993).

825 The EC number framework, which was developed
826 in 1992 by the “International Union of Biochemistry
827 and Molecular Biology”, is a contrasting approach
828 for functional annotation. Instead of focusing on
829 orthologous genes, EC numbers specify particular
830 enzyme-catalyzed reactions. An interesting charac-
831 teristic of this database is that these reactions can be
832 performed by non-homologous isofunctional enzymes
833 (Omelchenko et al. 2010). As of August 12th, 2020,
834 there were 6,520 EC numbers, which correspond to
835 one of four levels of granularity. For example, the ac-
836 cession EC 3.5.1.2 corresponds to glutaminases, while
837 the higher-level categories correspond to hydrolases
838 (3.-.-), that act on carbon-nitrogen bonds other than
839 peptide bonds (3.5.-), and that are in linear amides
840 (3.5.1.-). One major advantage of EC numbers is that
841 because they specify exact enzymatic reactions they
842 are straight-forward to link into pathway ontologies
843 based on reactions, such as MetaCyc pathways (Caspi
844 et al. 2013).

845 The Pfam database categorizes protein families,
846 which are protein regions that share sequence ho-
847 mology (Punta et al. 2012). Individual proteins
848 with multiple domains can thus belong to multiple
849 Pfam families. Each Pfam family is represented by
850 a hidden Markov model (HMM), which models the
851 likely amino acids at each residue and the likely
852 adjacent amino acids based on curated alignments
853 of representative protein sequence. This approach
854 identified homologous protein regions, which are
855 often hypothesized to have a shared evolutionary
856 history, but not necessarily. As of May 2020, there
857 were 18,259 Pfam families.

858 Lastly, TIGRFAMs are manually curated protein
859 families, which are also identified based on HMMs,
860 but also additional pertinent information (Haft et
861 al. 2003). As of September 16th, 2014, there were
862 4,488 TIGRFAMs. The distinguishing feature for
863 this database is that different information supple-
864 ments each HMM. For instance, certain TIGRFAM
865 are annotated based on species metabolic context
866 and neighbouring genes, while others are based on
867 validated functions from the scientific literature. This
868 database has been less commonly analyzed in recent
869 years and is best known as the annotation system
870 for early large-scale metagenomics projects (Venter
871 et al. 2004). Alternative approaches, such as the
872 FIGfam protein database are now more commonly
873 used than TIGRFAMs. FIGfams are based on a
874 similar approach, but instead of being manually
875 curated they are aggregated into isofunctional groups

876 based on shared roles in specific subsystems (Meyer
877 et al. 2009).

878 A recurrent question thus far has been that given
879 a range of comparable, or contrasting, bioinformatics
880 options, how is one to proceed? Fortunately, in the
881 case of selecting functional ontologies, the choice is
882 much clearer than other bioinformatics areas. Each
883 functional database typically excels for different pur-
884 poses. For instance, UniRef is useful for identifying
885 uncharacterized genes that may be of interest in
886 an environment, but quickly becomes challenging to
887 interpret and analyze in diverse communities.

888 In contrast, KEGG is useful for looking for shifts
889 in well-described functions at a high level, which
890 means this database is more robust to granular
891 functional diversity. Due to also being more robust
892 to granular functional diversity and because they
893 are more interpretable, pathway-level functions are
894 often of particular interest. For instance, obesity
895 is associated with an enrichment of phosphotrans-
896 ferase systems involved in carbohydrate processing
897 in human and mouse gut microbiomes (Turnbaugh
898 et al. 2008, 2009). This straight-forward explanation
899 quickly communicates the pertinent biological details,
900 which might be lost by focusing on less granular
901 levels.

902 However, it is worth noting that pathways identi-
903 fied based on DNA sequencing are merely theoretical
904 reconstruction based on the identified individual gene
905 families. Although there are several pathway recon-
906 struction approaches, they all require some mapping
907 from gene families or reactions to pathways. This
908 mapping can be structured, meaning that optional
909 and required contributors can be specified, or non-
910 structured, meaning that all genes and/or reactions
911 are treated equally.

912 The naïve approach for pathway reconstruction
913 is to assume that a pathway is present if any gene
914 or reaction involved is present in the community.
915 This was the predominant approach used for pathway
916 inference in early functional analyses (Moriya et al.
917 2007; Meyer et al. 2008) and in several pathway
918 inference tools such as PICRUSt (Langille et al.
919 2013). Pathway abundance under this framework
920 is calculated by summing the abundance of each
921 contributing gene family. This approach errs towards
922 avoiding missing the presence of a pathway, which is a
923 concern in metagenomes as key genes may be missing
924 due to mis-annotations. However, this approach
925 comes at the cost of spurious annotations. Based
926 on the naïve mapping approach the human genome
927 was previously annotated as including the KEGG
928 pathway equivalent of the reductive carboxylate cycle
929 (Ye and Doak 2011). This pathway is restricted to

930 autotrophic microbes and is similar to reversing the
931 Krebs cycle. Consequently, several gene families are
932 shared in both processes. Under the naïve mapping
933 approach, the presence of genes involved in the Krebs
934 cycle are also evidence for the predicted presence of
935 this atypical microbial pathway in humans. Similarly,
936 vitamin C biosynthesis would also be predicted in
937 humans based on the naïve approach (Ye and Doak
938 2011). However, the *GLO* gene, which encodes
939 the protein involved in the key last step of vitamin
940 C biosynthesis in mammals, is pseudogenized in
941 humans (Drouin et al. 2011), which makes vitamin
942 C biosynthesis impossible.

943 The Minimal set of Pathways (MinPath) ap-
944 proach is an approach developed to address this
945 issue (Ye and Doak 2011). This tool identifies
946 the smallest set of pathways, based on maximum
947 parsimony, that are required to explain the presence
948 of a set of proteins. In this way, the approach
949 is more conservative than naïve mapping and also
950 accounts for incomplete protein sets. This method
951 has been applied in numerous contexts, including for
952 the “HMP Unified Metabolic Analysis Network 2”
953 (HUMAN2) (Abubucker et al. 2012; Franzosa et
954 al. 2018) MGS gene family profiling and pathway
955 reconstruction framework. This popular framework
956 reconstructs pathways based on MinPath and infers
957 pathway abundance based on different approaches,
958 depending if the pathway mapping is structured. For
959 unstructured mappings, the arithmetic mean of the
960 upper half of individual gene family abundances is
961 taken to be the pathway abundance (Abubucker et
962 al. 2012). For structured mappings, the harmonic
963 mean of the key (i.e. required) genes families is
964 computed for pathway abundance (Franzosa et al.
965 2018). Both these approaches are motivated by the
966 need to be robust to variable abundance in alternative
967 gene families.

968 Although this approach for MGS pathway recon-
969 struction is commonly performed, it is important to
970 emphasize that it has not been universally accepted
971 and there remains disagreement about best-practices.
972 For example, “Evidence-based Metagenomic Path-
973 way Assignment using geNe Abundance DATA” (EM-
974 PANADA) is a method that addresses the same
975 issue as MinPath and HUMAN2 in a different way
976 (Manor and Borenstein 2017a). This method focuses
977 pathway reconstruction on distinguishing genes that
978 are shared with multiple pathways from those that
979 are unique to a single pathway. Pathway support
980 weightings are first given by the average abundance
981 of gene families unique to each given pathway. The
982 abundance of all shared gene families is then parti-
983 tioned between all pathways according to their rel-

984 ative support values. Pathway abundances are then
985 taken as the sum of the unique gene family relative
986 abundances and the partitioned relative abundances
987 of the shared gene families (Manor and Borenstein
988 2017a).

989 The exact reconstructed pathways and their re-
990 spective abundances differ depending on whether
991 naïve mapping, MinPath/HUMAN2, or EM-
992 PANADA are used. Validating pathway reconstruc-
993 tions is challenging without a gold-standard compar-
994 ison, particularly in metagenomes. Even in isolated
995 genomes, as demonstrated by the above examples of
996 the human pathway reconstructions, pathway recon-
997 struction is non-trivial. However, the advantage in
998 these cases is that experimental validation of pathway
999 reconstructions is possible (Francke et al. 2005;
1000 Oberhardt et al. 2008). Such validations would
1001 be possible if predictions are based on individual
1002 members of a microbiome, but it is less clear what
1003 experiments could validate pathways predicted for
1004 an overall community. In MGS data pathways are
1005 typically inferred as though all gene families were
1006 free to interact with each other. In other words,
1007 they are inferred as though there was universal cross-
1008 feeding. All three approaches described above are
1009 intended to be used for such community-wide gene
1010 family profiles. However, as mentioned above, this
1011 assumption is invalid because clearly not all proteins
1012 and metabolites in the microbiome can freely interact
1013 (McMahon 2015). The implications of this assump-
1014 tion being invalid remain unclear, but nonetheless it
1015 is an important caveat when interpreting pathway
1016 reconstruction data based on community-wide MGS
1017 data.

1018 This section would be incomplete without ad-
1019 dressing the most common discussion regarding mi-
1020 crobiome functional data: its ostensible high stability.
1021 Functional pathways are commonly at similar relative
1022 abundances across the same sample-types whereas
1023 taxonomic features, such as phyla, can substantially
1024 vary (Turnbaugh et al. 2009; Burke et al. 2011;
1025 HMP-consortium 2013; Louca et al. 2016). This
1026 functional consistency is often taken to be evidence
1027 of environmental selection for particular microbial
1028 functions (Turnbaugh et al. 2009; Louca and Doebeli
1029 2017). However, the validity of comparing variation
1030 between these two data types is rarely discussed.
1031 We and our colleagues investigated this question
1032 from a philosophical perspective and concluded that
1033 any meaningful comparison of the relative variation
1034 between taxonomic and functional profiles is likely
1035 impossible (Inkpen et al. 2017). This difficulty is
1036 largely because it is unclear which levels of gran-
1037 ularity would be meaningful to compare between

each data type. In other words, each data type is
qualitatively different from the other and the choice
of how to compare the two is based on arbitrary
decisions.

For instance, as discussed above, the sparsity
and number of possible functional categories differs
drastically across ontologies and sub-categories. We
demonstrated how observations of functional and
taxonomic stability are entirely dependent on how
function and taxa are defined (Inkpen et al. 2017).
We did this by comparing human stool sample
profiles at each possible taxonomic rank and also
each functional level for both the KEGG and UniRef
functional ontologies. As expected, phyla were less
stable across the samples than KEGG pathways,
but more stable than UniRef50 protein clusters.
However, this area remains an area of active debate.
Others have also argued that taxonomic variability
never unambiguously reflects functional variation,
which they believe is strong evidence for functional
conservation (Louca et al. 2018a). Nonetheless,
this example demonstrates once again the common
theme throughout this section: “function” has many
meanings.

Metagenome prediction methods

Ideally, analyses of microbial functions are based on
MGS data. However, predicted functions based on
16S rRNA gene (hereafter 16S) sequencing are often
analysed instead. **Metagenome** prediction, predicting
complete genomes for each individual ASV or taxon
weighted by their relative abundance, when based on
16S data is much cheaper than performing MGS.

There are additional advantages of predicted
metagenomes over actual MGS data. Namely, MGS
is often prohibitively expensive for samples where
host DNA overwhelms microbial DNA. The high
read depths required to yield sufficient microbial read
depths is infeasible in many cases (Gevers et al.
2014). Similarly, low-biomass samples are difficult
to accurately quantify with MGS, but they can
be profiled with PCR-based 16S sequencing. For
example, applying MGS to profile human tumours
is currently infeasible, but it is straight-forward to
apply 16S sequencing (Nejman et al. 2020). In
both cases, for host DNA contaminated and low-
biomass samples, metagenome prediction based on
16S profiles is a useful alternative to MGS.

However, metagenome prediction suffers from im-
portant drawbacks. The key problematic assumption
is that the marker gene used for predictions, typically

1089 the 16S, is strongly associated with genome content. 1143
1090 This broad assumption is correct: genera such as 1144
1091 *Lactobacillus* and *Desulfobacter* can be easily distin- 1145
1092 guished based on the 16S and they are enriched for 1146
1093 extremely different functions. Namely, *Lactobacillus* 1147
1094 can often perform lactic acid fermentation whereas 1148
1095 *Desulfobacter* can typically oxidize acetate to CO₂. 1149
1096 Such comparisons of characteristic functions between 1150
1097 distantly related taxa are uncontroversial. The 1151
1098 difficulty arises when approaches attempt to predict 1152
1099 entire genome contents for an entire community, 1153
1100 including for closely related taxa. 1154

1101 This issue is highlighted by classic DNA hy- 1155
1102 bridization experiments (Mandel 1966; Brenner 1156
1103 1973). These experiments were based on mixing 1157
1104 single-stranded DNA from two organisms and record- 1158
1105 ing the melting temperature required to separate the 1159
1106 strands. Higher melting temperatures are required to 1160
1107 break apart DNA that shares more complementary 1161
1108 bases connected by hydrogen bonds. Accordingly, 1162
1109 this approach provides a rough estimate of the genetic 1163
1110 distance between different strains or species. 1164

1111 An early comparison of these genetic distances 1165
1112 with 16S dissimilarity across 34 bacteria computed 1166
1113 a linear correlation of 0.728 (Devereux et al. 1990). 1167
1114 However, the relationship between these two metrics 1168
1115 is not linear: many bacteria with highly similar 1169
1116 16S genes have hybridization rates much lower than 1170
1117 70% (Stackebrandt and Goebel 1994), which is the 1171
1118 traditional cut-off for delineating species. This trend 1172
1119 has been corroborated across diverse prokaryotes 1173
1120 (Hauben et al. 1997, 1999; Kang et al. 2007). In 1174
1121 addition, a meta-analysis of 16S gene sequencing and 1175
1122 DNA hybridization data from 45 bacterial genera 1176
1123 further clarified these observations (Keswani and 1177
1124 Whitman 2001). This analysis established that 78% 1178
1125 of the variability in hybridization rates could be 1179
1126 accounted for by 16S similarity, based on a non- 1180
1127 linear model. However, they also identified that a 1181
1128 minority of hybridization rates were extremely poorly 1182
1129 predicted by 16S similarity (Keswani and Whitman 1183
1130 2001). 1184

1131 These observations agree well with genomic com- 1185
1132 parisons of strains, which can drastically differ in 1186
1133 genome content. For example, across 17 *E. coli* 1187
1134 genomes there are 13,000 genes that are variably 1188
1135 distributed and only 2,200 core genes (Rasko et 1189
1136 al. 2008). This enormous range of genomic vari- 1190
1137 ation is not reflected at the 16S level, where *E.* 1191
1138 *coli* strains are typically >99% identical (Suardana 1192
1139 2014). These genomic differences can translate to 1193
1140 enormous variation at higher taxonomic levels as well. 1194
1141 For instance, a comparison of the genomes from 11 1195
1142 *Yersinia* species found a range of genome sizes from 1196

3.7 - 4.8 megabases (Chen et al. 2010). A closer 1143
comparison of three pathogenic species of *Yersinia* 1144
determined that they shared 2,558 protein clusters 1145
while 2,603 were variably distributed. These species- 1146
level differences are also not proportionally reflected 1147
by divergence in *Yersinia* species 16S genes, which 1148
are typically >97% identical (Ibrahim et al. 1993). 1149
These examples highlight that 16S similarity can be 1150
a poor predictor of genomic similarity. This issue 1151
is compounded when there are divergent 16S copies 1152
within the same genome, although typically these are 1153
>99.5% identical (Větrovský and Baldrian 2013). 1154

1155 Variation in gene content within a taxonomic 1156
lineage is a recurrent observation across microbial 1157
communities. Variably present genes are often linked 1158
to putative niche-specific adaptations (Wilson et al. 1159
2005), such as genes affecting antibiotic resistance 1160
(Kallonen et al. 2017), carbohydrate catabolism 1161
(Arboleya et al. 2018), and wound healing (Kalan 1162
et al. 2019). Based on these and other observations, 1163
the understanding of bacterial genomic content has 1164
shifted from that of a static genome to a pan-genome, 1165
consisting of core and variable genes (Tettelin et 1166
al. 2005). Variably present genes are transmitted 1167
between genomes through horizontal gene transfer, 1168
which typically occurs between closely related organ- 1169
isms (Popa and Dagan 2011). However, horizontal 1170
gene transfer can also occur between distantly related 1171
organisms, such as between different bacterial phyla 1172
(Beiko et al. 2005; Kloesges et al. 2011; Martiny et 1173
al. 2013). 1174

1175 The high variability between bacterial genomes 1176
and extensive horizontal gene transfer highlights 1177
the major challenges facing metagenome predic- 1178
tion. Despite these challenges, interest in performing 1179
metagenome predictions has continued, supported 1180
by several observations. First, although there are 1181
important outliers, 16S sequence identity does log- 1182
arithmically correlate well with the average nu- 1183
cleotide identity between genomes, with an R² of 1184
0.79 (Konstantinidis and Tiedje 2005). Second, 16S 1185
sequence similarity does provide some information 1186
on the ecological similarity of bacteria (Chaffron et 1187
al. 2010). This was demonstrated by the fact that 1188
co-occurring environmental bacteria are more likely 1189
to have similar 16S sequences. In addition, overall 1190
differences in inferred KEGG pathway potential are 1191
strongly associated with 16S divergence (Chaffron et 1192
al. 2010). Last, within a given environment, such 1193
as the human gut, 16S divergence was shown to be 1194
particularly predictive of divergence in average gene 1195
content (Zaneveld et al. 2010). 1196

1197 Originally, metagenome prediction workflows 1198
were based on matching 16S sequences to reference 1199

1197 genomes. By taking the best matching genome or
1198 averaging across genomes with similar sequences, a
1199 predicted genome annotation can be acquired for all
1200 16S sequences (Figure 1d). To infer the metagenome
1201 profile one must simply multiply the predicted
1202 genome annotations for each 16S sequence by the
1203 abundance of each 16S sequence in the metagenome.
1204 In addition to predicting microbial functions linked to
1205 Crohn’s disease (Morgan et al. 2012), this approach
1206 has also been used to profile diet-related microbial
1207 functions across mammals (Muegge et al. 2011)
1208 and the functions of invasive bacteria within corals
1209 (Barott et al. 2012). Although bioinformatics tools
1210 for metagenome prediction are now typically used for
1211 performing this task, this 16S-matching approach is
1212 still used for custom analyses (Verster and Borenstein
1213 2018; Bradley and Pollard 2020).

1214 The first metagenome prediction tool to expand
1215 beyond this approach, and specifically intended for
1216 16S sequencing data, was “Phylogenetic Investigation
1217 of Communities by Reconstruction of Unobserved
1218 States” (PICRUSt1) (Langille et al. 2013). This
1219 tool is based on leveraging classical ancestral-state
1220 reconstruction methods, which have been widely used
1221 in phylogenetics (Zaneveld and Thurber 2014). The
1222 crucial extension of this framework is to extend
1223 trait predictions from internal, or ancestral, nodes
1224 in a phylogenetic tree to tips with unknown trait
1225 values. This approach has been termed hidden-
1226 state prediction (HSP) (Zaneveld and Thurber 2014).
1227 We recently published a major update to PICRUSt,
1228 called PICRUSt2 (Douglas et al. 2020). The key
1229 improvement in PICRUSt2 is that predictions can
1230 be made for novel 16S sequences with this tool
1231 and custom databases can be more easily used for
1232 analyses.

1233 PICRUSt1 introduced the step of normalizing
1234 relative abundances by the predicted number of 16S
1235 copies within each genome, which is intended to
1236 control biases in 16S sequencing due to copy number
1237 (Farrelly et al. 1995). Importantly, although 16S
1238 copy number correction has become a common step
1239 for metagenome prediction (Angly et al. 2014),
1240 accurately predicting 16S copy number is particularly
1241 challenging. An independent validation of several
1242 16S copy number prediction methods, including PI-
1243 CRUSt1, identified poor agreement of predicted copy
1244 numbers against existing reference genomes (Louca
1245 et al. 2018b). In some cases, less than 10% of the
1246 variance in actual 16S copy number was explained
1247 by these predictions. In addition, these predictions
1248 were often only slightly correlated between prediction
1249 methods.

1250 Since PICRUSt1 was published a number of

1251 similar metagenome prediction tools have been de-
1252 veloped. All of these approaches aim to capture
1253 the shared phylogenetic signal in the distribution of
1254 functions across taxa. These tools include: PanFP
1255 (Jun et al. 2015), Piphillin (Iwai et al. 2016; Narayan
1256 et al. 2020), PAPRICA (Bowman and Ducklow
1257 2015), and Tax4Fun2 (Wemheuer et al. 2020).

1258 These metagenome prediction tools have primar-
1259 ily been validated by comparing how well the pre-
1260 dicted gene family abundances they output correlate
1261 with the abundances of gene families identified in
1262 MGS data from the same samples. This approach
1263 generally identifies high correlations between the two
1264 profiles. For example, predicted KOs output by
1265 PICRUSt1 based on Human Microbiome Project
1266 (HMP) samples were highly correlated with the
1267 matching MGS-identified data (Spearman $r = 0.82$)
1268 (Langille et al. 2013). Importantly, a high Spearman
1269 correlation is actually expected by chance in these
1270 comparisons simply because many genes are common
1271 in most environments while others are usually absent
1272 or rare. Upon comparing to this expectation the
1273 predictions are still significantly better than expected
1274 by chance, but only slightly (Douglas et al. 2020).
1275 Nonetheless, based on this approach, we found that
1276 PICRUSt2 performed marginally better than other
1277 tools (Douglas et al. 2020). However, it is noteworthy
1278 that Piphillin, which represents a much simpler
1279 approach based on a nearest-neighbour approach,
1280 performed only slightly worse overall and better in
1281 some contexts.

1282 An alternative approach for evaluating these
1283 methods is based on the concordance of differen-
1284 tial abundance results between actual and predicted
1285 metagenomics profiles. When we conducted this anal-
1286 ysis while validating PICRUSt2, we found that dif-
1287 ferential abundances tests on metagenome prediction
1288 tools agreed only moderately well with matching tests
1289 based on actual MGS data (Douglas et al. 2020).
1290 This is a crucial point to appreciate when analyzing
1291 metagenome prediction data; even though the overall
1292 predicted profiles might correlate with MGS profiles,
1293 the results from differential abundance testing might
1294 nonetheless be quite different. We also observed
1295 high variation across datasets in concordance between
1296 MGS and 16S-based predictions. In other words,
1297 differential abundance testing on predicted profiles
1298 resulted in fair agreement with MGS data on some
1299 datasets while disagreeing almost entirely on others.
1300 In addition, researchers performing independent work
1301 in this area have identified conflicting signals of how
1302 well individual metagenome prediction tools perform
1303 (Narayan et al. 2020; Sun et al. 2020). These
1304 observations might again reflect the high variation

1305 across datasets in how well prediction profiles agree
1306 with MGS results.

1307 **Current state of the integration** 1308 **of taxonomic and functional** 1309 **data types**

1310 The above discussion has described the many faces
1311 of microbiome data types. Taxonomic and functional
1312 microbiome data are typically generated indepen-
1313 dently, but in some cases can be directly linked.
1314 Regardless of the exact processing workflow for these
1315 data types, we have yet to address one question: how
1316 are they integrated?

1317 For independent taxonomic and functional data
1318 types this is largely done anecdotally. For example,
1319 this is commonly done in regards to the nine genera
1320 that are the primary producers of short-chain fatty
1321 acids (SCFAs) in the human gut (Moya and Ferrer
1322 2016). SCFA levels have long had an ambiguous
1323 link with Crohn’s disease (CD) (Treem et al. 1994),
1324 although they are typically negatively associated with
1325 disease activity (Venegas et al. 2019). Due to this
1326 association, there has been long-standing interest in
1327 identifying microbial taxa that are associated with
1328 altered SCFA levels. Accordingly, CD microbiome
1329 studies commonly hypothesize that shifts in the re-
1330 lative abundance of any known SCFA-producing taxa
1331 likely cause altered SCFA levels. For example, *Fae-*
1332 *calibacterium prausnitzii* is a well-known commensal
1333 SCFA-producer in the human gut and is consistently
1334 found at lower levels in the CD patient microbiomes
1335 (Wright et al. 2015). Although potential links
1336 between lower levels of this species, in addition to
1337 other taxa such as *Roseburia* (Laserna-Mendieta et
1338 al. 2018), and SCFA levels are often discussed, this
1339 is rarely formally investigated.

1340 More often, anecdotal links between function and
1341 taxa are based on observed associations between sig-
1342 nificant features. Several such cases have previously
1343 been noted as representative examples (Manor and
1344 Borenstein 2017b). For instance, *Propionibacterium*
1345 *acnes* has been identified as strongly correlated with
1346 NADH dehydrogenase levels in the skin microbiome
1347 (Oh et al. 2014). Consequently, this species was
1348 implicated as the likely cause for changes in NADH
1349 dehydrogenase levels. Similarly, *Bacteroides thetaio-*
1350 *taomicron* relative abundance has been identified as
1351 positively correlated with microbial genes involved
1352 with the degradation of complex sugars and starch
1353 in the infant gut (Bäckhed et al. 2015). Based
1354 on this observation, this species was hypothesized

1355 to be the key contributor to increased levels of
1356 these degradation genes. Such insights are valuable,
1357 but as previously discussed (Manor and Borenstein
1358 2017b), these anecdotal links alone are not convincing
1359 evidence that particular taxa are the primary contrib-
1360 utors to functional shifts.

1361 Linked taxonomic and functional data alone is not
1362 sufficient to resolve this issue. There are substantial
1363 challenges facing the integration of these data types
1364 besides simply generating a combined format. For
1365 example, two massive datasets have recently been
1366 published as part of the next iteration of the Human
1367 Microbiome Project. Both datasets include numerous
1368 sequencing and profiling technologies, including 16S
1369 and MGS, from the stool and various body-sites
1370 of **IBD** (Lloyd-Price et al. 2019) and individuals
1371 with pre-diabetes (Zhou et al. 2019). However, in
1372 each case there was little integration of microbiome
1373 functional and taxonomic data types. Instead, these
1374 features were largely tested independently, despite
1375 the availability of links between the data types,
1376 and associations between top features were discussed
1377 (Lloyd-Price et al. 2019; Zhou et al. 2019).

1378 In contrast to these examples, there have been
1379 calls for improved integration of these microbiome
1380 data types, which is rooted in a systems-level biology
1381 outlook (Greenblum et al. 2013). “Functional
1382 Shifts’ Taxonomic Contributors” (FishTaco) is one
1383 bioinformatics method developed for this purpose,
1384 which quantifies taxonomic contributions to func-
1385 tional shifts (Manor and Borenstein 2017b). One
1386 major application of this approach is to distinguish
1387 two explanations for why a function might be at
1388 high relative abundance (Figure 2). First, a function
1389 might be higher in relative abundance simply because
1390 it hitchhiked on the genome of a taxon that bloomed
1391 for other reasons. In contrast, an alternative explana-
1392 tion might be that many taxa performing the same
1393 function gained a growth advantage and thus grew
1394 in relative abundance. FishTaco can also identify
1395 functions that have grown in relative abundance
1396 simply because microbes that do not encode it are
1397 at lower levels.

1398 FishTaco works by first identifying significant
1399 shifts in functional abundances with a standard
1400 differential abundance test, typically a Wilcoxon test.
1401 Subsequently, a permutation analysis is undertaken,
1402 which consists of randomly shifting the relative abun-
1403 dance of a subset of taxa, while maintaining the
1404 rest. A large collection of such permutations is
1405 performed, which include permutations of single and
1406 multiple taxa in different replicates. Based on this
1407 approach an estimate of the relative contribution of
1408 each taxon to a functional shift can be estimated

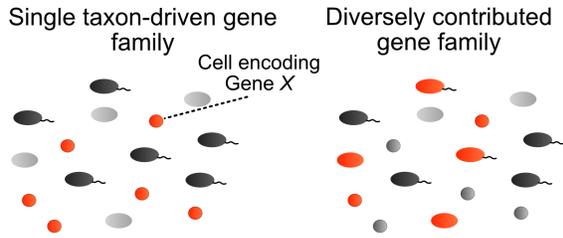


Figure 2: Two explanations for why a gene family might be at higher relative abundance that would be impossible to distinguish without joint taxonomic and functional data. Microbes encoding the gene of interest (Gene X) are indicated in red. This diagram contrasts how a gene family might be blooming due to a single taxon (left) versus a diverse set of taxa (right). The importance of distinguishing these scenarios is underappreciated: in the second case it is more likely the gene family itself that confers a growth or survival advantage in the environment. Note however that these are not the only two reasons why the relative abundance of a gene family might be at high levels in an environment.

(Manor and Borenstein 2017b). These relative contributions are then presented as stacked bar charts breaking down the direction and magnitude of each functional contribution. These visualizations help distinguish when a functional shift is due to the enrichment or depletion of taxa and also which sample grouping the shift occurred within. This approach was motivated by Shapley values, which were introduced in game-theory to summarize the contribution of each player in a multiplayer game (Shapley 1953). Specifically, FishTaco leverages a modified version of this approach that enables the contribution of individual features to be estimated in large datasets without exhaustively testing every possible permutation (Keinan et al. 2004).

FishTaco represents an important advancement in integration and improved interpretability of taxonomic and functional microbiome data. However, it nonetheless suffers from major limitations. First, although the taxonomic breakdown of contributors to a function is valuable, the FishTaco approach requires significant functions to be identified based on the relative abundance of individual gene families and pathways. This is done by systematically testing all functions across the entire metagenome, which is problematic when performed with a non-compositional approach like a Wilcoxon test. This approach also treats gene families under the bag-of-genes model, which is inappropriate, as discussed above. An improved method would conduct a com-

positionally sound analysis and integrate taxonomic information when identifying significant functions.

An alternative method is phylogenize, which does address each of these issues (Bradley et al. 2018; Bradley and Pollard 2020). This approach tests for significant associations between the presence of a taxa within a given sample grouping and the probability that a taxon encodes a given gene family. This is performed through phylogenetic linear regression, which accounts for the genetic similarity of co-occurring taxa that might trivially be due to a shared evolutionary history. A separate phylogenetic linear model is fitted for each gene family. The key distinction of this approach from a normal linear model is that instead of the residuals being independent and normally distributed, they covary so that phylogenetically similar microbes have higher covariance (Bradley et al. 2018). This overall approach was partially motivated by an attempt to address a similar problem by comparing the species and gene trees of gut and non-gut microbes (Lozupone et al. 2008). Based on simulated random data (i.e. data with no real functional shifts) the phylogenize authors demonstrated that performing standard linear models without controlling for phylogenetic structure results in false positive rates ranging from 20% - 68%. In contrast, controlling for phylogenetic structure with phylogenize resulted in a uniform P-value distribution and an appropriate false positive rate of 5%. One interesting feature is that phylogenize does not directly analyze relative abundances. Instead, the tool converts taxa relative abundance into one of three formats: (1) binary presence/absence across all samples, (2) overall prevalence within each sample grouping, (3) or the specificity within each sample grouping (Bradley et al. 2018).

Although phylogenize is undeniably an invaluable contribution to microbiome data analysis, it also has several limitations. First, information on taxa abundance is discarded entirely in favour of presence/absence data. From one perspective this is an advantage; eliminating taxa relative abundances enables phylogenize to circumvent compositionality issues. However, relative abundance data is often more important to investigate, because key taxonomic shifts might not be detected by presence/absence alone. In addition, phylogenize reports significant gene families for each phylum in a dataset. This is performed to reduce the memory usage and to enable phylum-specific rates of evolution for each function (Bradley et al. 2018). This focus on the phylum level makes the results difficult to interpret for two reasons. First, it is insufficiently broad, because it limits the potential to identify functions distributed

1493 across multiple phyla that might be linked with a
1494 condition of interest. From another perspective,
1495 this focus on the phylum level is also not specific
1496 enough; although phylum-function associations are
1497 valuable they do not provide information on the
1498 relative contributions of lower-level taxa, such as
1499 species, to the association. Accordingly, there is room
1500 for improvement in both the statistical analysis and
1501 interpretation of the phylogenize approach.

1502 Despite the availability of approaches for integrat-
1503 ing functional and taxonomic data, they have yet to
1504 become a mainstay of microbiome analyses. However,
1505 it is becoming common to visualize stacked bar-charts
1506 of taxonomic contributors to functions of interest.
1507 This is typically performed on predicted metagenome
1508 output by PICRUSt or alternatively on HUMAnN2
1509 output, although this could be performed with any
1510 linked taxa-function data. As discussed above, the
1511 HUMAnN2 pipeline includes a step for identifying
1512 particular strains in MGS dataset, which allows gene
1513 families to be linked to those strains (Franzosa et al.
1514 2018). In some cases this approach enables complete
1515 links between taxa and function to be identified. For
1516 instance, *F. prausnitzii* was shown to be the obvious
1517 principal contributor to glutaryl-CoA biosynthesis
1518 in the HMP gut MGS samples (Franzosa et al.
1519 2018). However, more commonly there are numerous
1520 taxonomic contributors to a single given function,
1521 and it is difficult to interpret which taxa are the
1522 key contributors by looking at visualizations alone.
1523 Nonetheless, even in the presence of many taxonomic
1524 contributors, the HUMAnN2 authors demonstrated
1525 that these visualizations can provide information
1526 about the diversity of taxa contributing to a function,
1527 termed the contributinal diversity (Franzosa et al.
1528 2018). This is most often quantified with the Gini-
1529 Simpson index, which is the complement of Simpson's
1530 evenness (Jost 2006).

1531 Contributinal diversity has been shown to be a
1532 useful approach for delineating housekeeping path-
1533 ways encoded by many taxa, intermediate pathways,
1534 and those rarely encoded, which can correspond
1535 to opportunists or keystone species. For instance,
1536 *F. prausnitzii* has previously been linked with sev-
1537 eral human microbiome pathways identified through
1538 MGS that have intermediate contributinal diversities
1539 (Abu-Ali et al. 2018). When present, this species
1540 tended to contribute the majority of all pathways it
1541 encoded.

1542 This approach has also been valuable for profiling
1543 shifts in the contributions to microbial pathways
1544 over time, such as in the infant gut profiled with
1545 MGS (Vatanen et al. 2018). In this case, several
1546 microbial pathways, such as siderophore biosynthe-

1547 sis, were found to display decreasing contributinal
1548 diversity with age. This is an interesting observation
1549 because siderophores are costly to produce but are
1550 highly beneficial in the human gut. In particular,
1551 siderophores can confer a strong benefit to multi-
1552 ple community members, including those that do
1553 not produce siderophores, by providing access to
1554 iron. Siderophores have previously been presented as
1555 microbial functions whose distribution is consistent
1556 with the Black Queen Hypothesis (Morris et al.
1557 2012). This hypothesis states that adaptive gene loss
1558 may occur for functions that are costly to produce,
1559 provided that the function is provided by other
1560 community members. This hypothesis was discussed
1561 in the context of the infant microbiome as an expla-
1562 nation for why siderophore contributinal diversity
1563 decreases over time (Vatanen et al. 2018): perhaps
1564 gene loss confers an adaptive benefit by avoiding the
1565 production of a costly metabolite. Although this is
1566 an interesting hypothesis, a less controversial inter-
1567 pretation of this result is simply that siderophores
1568 became less stably encoded over time in the profiled
1569 samples.

1570 Related to this point, two additional metrics
1571 have also been developed to summarize the stability
1572 of taxonomic contributions to microbial functions
1573 (Eng and Borenstein 2018). More specifically, these
1574 metrics are intended to summarize functional robust-
1575 ness across samples, which is the stability in the
1576 relative abundance for a given function in response
1577 to taxonomic perturbation. This is performed by
1578 generating a taxa-response curve that describes the
1579 average change in functional relative abundances in
1580 response to taxonomic perturbations of different mag-
1581 nitudes. Two metrics are then computed based upon
1582 these curves: attenuation and buffering. Attenuation
1583 captures how rapidly a function shifts with increasing
1584 taxonomic perturbation magnitudes. In contrast,
1585 buffering represents how well functional shifts are
1586 suppressed at smaller taxonomic perturbation mag-
1587 nitudes.

1588 Applying these metrics to PICRUSt-predicted
1589 metagenomes from 16S sequencing of human body
1590 sites, validated by a subset of MGS samples, yielded
1591 several novel perspectives. First, attenuation and
1592 buffering were conserved across body sites for micro-
1593 bial house-keeping pathways but varied for several
1594 others. For instance, robustness in the biosynthesis
1595 of unsaturated fatty acids varied substantially across
1596 body sites. In addition, human gut samples were
1597 found to have higher values of both attenuation and
1598 buffering than compared to vaginal samples. These
1599 trends were shown to be driven by more than simply
1600 lower richness in vaginal samples by subsampling

1601 to comparable diversity levels across each body-site
1602 (Eng and Borenstein 2018). These observations are
1603 consistent with the controversial hypothesis that mi-
1604 crobial communities may be under varying selection
1605 strengths for functional robustness, depending on the
1606 environment (Naeem et al. 1998; Ley et al. 2006).

1607 The development of these metrics for summariz-
1608 ing functional contributions represent an important
1609 goal of microbiome research, which is to leverage
1610 sequencing data to yield novel biological insights. In
1611 contrast, another major goal is to answer a more
1612 practical question: how useful is microbiome data for
1613 classification and prediction tasks?

1614 There is great interest in applying machine
1615 learning approaches to microbiome sequencing data
1616 (Knights et al. 2011). Most commonly this is
1617 performed with either Support Vector Machine or
1618 Random Forest (Breiman 2001) models. Applications
1619 of these and other machine learning approaches to
1620 microbiome data are primarily aimed at distinguish-
1621 ing samples from different environments or disease
1622 states (Zhou and Gallins 2019). Taxonomic features
1623 are the focus of most such microbiome-based machine
1624 learning approaches, which is true for both 16S
1625 (Duvall et al. 2017) and MGS (Pasolli et al.
1626 2016) data. However, on a growing number of
1627 occasions machine learning is focused on functional
1628 data types. For example, a recent MGS meta-analysis
1629 identified informative functional biomarkers across
1630 several human diseases by applying machine learning
1631 approaches to functional data types (Armour et al.
1632 2019). Regardless of the data type, models trained on
1633 microbiome data typically have low generalizability
1634 across independent cohorts (Sze and Schloss 2016;
1635 Douglas et al. 2018), although there are exceptions.

1636 One major exception is microbiome-based mod-
1637 elling of colorectal cancer, which in one investigation
1638 was shown to be generalizable across five independent
1639 datasets (Wirbel et al. 2019). This landmark study
1640 also systematically compared the utility of functional
1641 and taxonomic data types in these models and found
1642 them to be comparable overall. This finding is
1643 consistent with a past comparison of the classifica-
1644 tion performance of 16S-based taxa and predicted
1645 metagenome data (Ning and Beiko 2015). In the case
1646 of predicted metagenomes, which are based on 16S
1647 profiles, it is perhaps less surprising that they yield
1648 comparable classification performance. However,
1649 with MGS data in particular it might be possible to
1650 detect robust, informative functions that might be
1651 undetectable with taxonomy alone due to taxonomic
1652 variability (Doolittle and Booth 2017).

1653 Despite this great interest in applying machine
1654 learning to different microbiome data types, there

1655 has been little focus on integrating across them.
1656 The aforementioned comparison of 16S-based taxa
1657 and predicted functions is one exception where a
1658 hybrid classification model of both data types was
1659 created (Ning and Beiko 2015). In this case, there
1660 was a small increase in classification performance
1661 for distinguishing nine human oral sub-locations.
1662 The original OTU and KO-based models yielded
1663 accuracies of 76.2% and 76.1%, respectively, while
1664 the hybrid model resulted in an accuracy of 77.7%
1665 (Ning and Beiko 2015). This result indicates that
1666 predicted functions may provide some additional
1667 information in combination with taxonomic data, but
1668 the consistency and biological significance of this
1669 small effect remains unclear. Further investigation
1670 into the integration of these data types within a
1671 machine learning context is needed to ensure that the
1672 highest-quality models possible are constructed.

1673 Outlook

1674 Herein we have described the unique characteristics
1675 of microbiome DNA data types and many of the ap-
1676 proaches that have been proposed for their analysis.
1677 Throughout we have emphasized two ideas. First,
1678 increased integration of taxonomic and functional
1679 microbiome data types is needed. And second,
1680 there is often high variation in the results between
1681 microbiome data analysis pipelines.

1682 Regarding the first point, we believe that several
1683 of the tools described above, such as FishTaco and
1684 phylogenize, largely solve the issue of how to jointly
1685 investigate taxa and functions. Increased usage and
1686 development of these and other related tools would
1687 greatly help with the interpretability of microbiome
1688 data.

1689 One area where further development is particu-
1690 larly needed is in the context of classification models,
1691 where little work has been conducted to systemat-
1692 ically link taxa and functions appropriately. One
1693 exception was a classification approach based on gene
1694 families that identified predictive genes and then sub-
1695 sequently identified metagenome assembled genomes
1696 within a given dataset enriched for these genes (Rah-
1697 man et al. 2018). However, this approach still relied
1698 on follow-up analyses rather than integrating the data
1699 types. Instead, an improved approach could be based
1700 on explicitly leveraging the hierarchical nature of
1701 microbiome data types. This is because functional
1702 and taxonomic data types independently form clear
1703 hierarchical structures (e.g. Pathway - Gene and
1704 Phylum - Class - Order, etc.). The connection
1705 between taxa and gene families and pathways is

1706 more complex, but nonetheless, links between groups
1707 of strains or ASVs and microbial functions can be
1708 defined. A modified machine learning framework
1709 that explicitly accounted for these relationships could
1710 result in more interpretable outputs.

1711 Regardless of the specific tool, microbiome re-
1712 searchers should move towards more integration of
1713 taxonomic and functional data. It is odd to distin-
1714 guish between functional and taxonomic datatypes in
1715 the first place: they are inextricably linked after all.
1716 The term “metagenome” itself is in some ways unfor-
1717 tunate as it implies that the genetic information for
1718 all organisms in a community can be simultaneously
1719 analyzed in a coherent way, without partitioning
1720 genes into genomes. This may be valid for high-level
1721 pathways but for generating hypotheses regarding
1722 specific gene families it is too often misleading. This
1723 perspective is becoming more common, as the avail-
1724 ability of metagenome-assembled genomes increases
1725 (Frioux et al. 2020).

1726 The other common thread throughout this
1727 manuscript has been that technical variation in mi-
1728 crobiome data analyses means that making robust
1729 biological inferences, especially regarding specific mi-
1730 crobial features, is challenging. Indeed, the lack
1731 of standardization in microbiome data analysis has
1732 previously been strongly criticized. An assessment
1733 of numerous papers attempting to define standard
1734 pipelines concluded that there was disturbingly little
1735 consensus (Pollock et al. 2018). This is true for
1736 many steps related to the processing, sequencing,
1737 and analysis of microbiome data. For instance,
1738 there have been contradictory results regarding the
1739 efficacy of different extraction protocols (Salonen et
1740 al. 2010). In particular, underrepresentation of
1741 Gram-positives has been observed (Maukonen et al.
1742 2012), which may be partially resolved by using
1743 bead-beating extraction protocols (Guo and Zhang
1744 2013). There is also substantial technical variation
1745 related to bioinformatics choices, which represent the
1746 final steps of a microbiome project. For example,
1747 as discussed above, the bioinformatics choices made
1748 when performing differential abundance testing on
1749 microbiome data can have severe impacts on any
1750 interpretations (Thorsen et al. 2016; Hawinkel et al.
1751 2019).

1752 We have encountered similar issues with our
1753 work, most strikingly when investigating pediatric
1754 Crohn’s disease patients’ microbiome profiles (Dou-
1755 glas et al. 2018). An important characteristic of
1756 these data was that 98% of the sequenced reads
1757 mapped to the human genome. This characteristic
1758 made taxonomic profiling of these data especially
1759 prone to false positives. In particular, an initial

1760 draft of our manuscript was based on profiles that
1761 included large proportions of viral-identified DNA
1762 and matches to certain eukaryotic parasites. We were
1763 initially excited about these observations, because
1764 the abundances of these non-prokaryotic taxa were
1765 discriminative for classifying patient disease state
1766 and treatment response. However, the exact taxa
1767 identified were peculiar: they were predominately
1768 represented by a range of plant-associated viruses
1769 and the eukaryotic genus *Plasmodium*, which is best
1770 known as including the causative agent for malaria,
1771 *Plasmodium falciparum*. Upon closer investigation it
1772 became clear that this signal was driven entirely by
1773 a difference in how reads were mapped to lineage-
1774 specific marker genes. Altering the parameter choice
1775 from local to global mapping entirely removed these
1776 taxa. This relatively small difference in parameter
1777 choice appeared to only affect our data and not
1778 more typical microbiome datasets, which we believe
1779 was due to the high proportion of human DNA
1780 in our data. Although this error was moderately
1781 embarrassing, it was more importantly an example
1782 of how easily a single parameter setting can result
1783 in starkly different biological interpretations. In this
1784 case the difference was driven by an option used for
1785 a single bioinformatics tool.

1786 Such inconsistencies in microbiome analyses have
1787 previously been identified and been shown to make
1788 meaningful comparisons across studies challenging.
1789 For instance, associations between obesity and the
1790 human microbiome are commonly discussed as sup-
1791 port for the utility of considering microbial links
1792 with human disease, despite inconsistencies across
1793 studies (Castaner et al. 2018; Muscogiuri et al.
1794 2019). These inconsistencies are typically explained
1795 due to confounding variables that may differ between
1796 patient cohorts. Although this is a valid explanation,
1797 it is likely that technical variation, including in
1798 terms of bioinformatics analyses, also drives these
1799 inconsistencies. For instance, a meta-analysis of ten
1800 obesity human microbiome datasets identified only
1801 extremely weak signals when re-analyzing all datasets
1802 with a standardized approach (Sze and Schloss 2016).
1803 This finding greatly contrasts with how these studies
1804 were originally presented and again highlights how
1805 variation in bioinformatics can greatly affect how to
1806 biologically interpret microbiome data.

1807 Similarly lower alpha diversity in stool micro-
1808 biomes has been frequently linked with disease states
1809 (Mosca et al. 2016). These observations are intu-
1810 itively reasonable as reduced alpha diversity could
1811 enable pathogens to bloom (Vincent et al. 2013) or
1812 represent differences in resource availability (Turn-
1813 baugh et al. 2009). However a re-analysis of data

1814 from 28 studies representing ten diseases was unable
1815 to identify evidence for links between alpha diversity
1816 and disease states (Duvallet et al. 2017). The
1817 exceptions were diarrheal diseases and inflammatory
1818 bowel diseases.

1819 Such inconsistencies across analyses on the same
1820 data are gradually coming to the forefront of the
1821 microbiome field (Allaband et al. 2019). Indeed,
1822 a recent plea for improved standardization has been
1823 made to enable better comparisons across studies
1824 (Hill 2020). This is a commendable goal, but given
1825 the diversity of opinions regarding best-practices
1826 (Callahan et al. 2016b; Knight et al. 2018; Schloss
1827 2020), it is difficult to coherently recommend a single
1828 workflow for analyses at the moment. Accordingly,
1829 further work and benchmarking of different bioin-
1830 formatics is needed to convincingly argue for best
1831 practices in microbiome data analysis.

1832 Until a clear consensus is reached it is the res-
1833 sponsibility of microbiome researchers to make the
1834 caveats and challenges facing this area clear to read-
1835 ers and newcomers to the field. This is crucial given
1836 the widespread interest in studying microbiomes
1837 through DNA sequencing; the number of microbiome
1838 sequencing-related publications continues to rapidly
1839 grow. This is in tandem with funding for these
1840 projects, which has steadily increased in the USA
1841 from at least 2007 to 2016 (NIH 2019). According to
1842 the US National Health Institute, there was US\$766
1843 million dollars invested in microbiome research in
1844 2019, which was the 63rd most highly funded health-
1845 related research category out of 291. Although
1846 comparing across research categories of varying gran-
1847 ularity is difficult, it is noteworthy that microbiome
1848 research was more highly funded than both breast
1849 cancer and Alzheimer’s disease research. Import-
1850 antly, an increased interest in microbiome research is
1851 warranted: recent technological developments are en-
1852 abling improved investigations into microbial biology.
1853 However, as the monetary investment and research
1854 hours dedicated to microbiome research grows, it is
1855 crucial that scientists ensure the best use of these
1856 resources. Open discussions on the many contentious
1857 aspects of microbiome data analysis would help with
1858 this issue. Indeed, such clarifications by leaders in the
1859 microbiome field are starting become more common
1860 (Allaband et al. 2019). However, although these
1861 contributions are valuable, they do not adequately
1862 address the problem. In particular, instead of men-
1863 tioning these issues in passing, inconsistencies be-
1864 tween bioinformatics workflows should be emphasized
1865 more clearly for the benefit of the uninitiated.

1866 Another practical improvement would be to nor-
1867 malize, and potentially require, explicit summaries

1868 of the effects of technical variation on any biological
1869 interpretations reported in microbiome studies. This
1870 is impossible to capture entirely, but it could be
1871 done by comparing how key results change depending
1872 on a subset of representative bioinformatics choices.

1873 For instance, researchers could compare how insights
1874 change depending on the combinations of denoising
1875 tools and differential abundance methods that they
1876 have applied when analyzing 16S data. Although
1877 these changes would result in increased workloads
1878 when conducting analyses and when communicating
1879 results, they would help ensure that any major bio-
1880 logical findings are at least robust to a representative
1881 set of bioinformatics choices.

1882 Regardless of which approach is taken to address
1883 these issues, the most important point is that action
1884 is needed on this front. The variation between bioin-
1885 formatics methods is undeniable and unfortunately
1886 reflects a reproducibility crisis facing microbiome
1887 data analysis.

1888 Acknowledgements

1889 We would like to thank the following individuals for
1890 feedback on sections of this manuscript: Dr. Robert
1891 Beiko, Dr. Zhenyu Cheng, Dr. André Comeau,
1892 Casey Jones, Jacob Nearing, Dr. Laura Parfrey,
1893 Dr. Andrew Stadnyk, Chris Tang, and Dr. Robyn
1894 Wright.

1895 Conflicts of interest

1896 We declare that we have no conflicts of interest with
1897 the content of this article.

1898 References

- 1899 Abu-Ali GS et al. 2018. Metatranscriptome of
1900 human faecal microbial communities in a cohort
1901 of adult men. *Nat. Microbiol.* 3:356-366.
- 1902 Abubucker S et al. 2012. Metabolic reconstruction
1903 for metagenomic data and its application to
1904 the human microbiome. *PLOS Comput. Biol.*
1905 8:e1002358.
- 1906 Aitchison J. 1982. The Statistical Analysis of Com-
1907 positional Data. *J. R. Stat. Soc. Ser. B.* 44:139-
1908 177.
- 1909 Allaband C et al. 2019. Microbiome 101: Studying,
1910 Analyzing, and Interpreting Gut Microbiome
1911 Data for Clinicians. *Clin. Gastroenterol. Hepa-
1912 tol.* 17:218-230.

1913 Ambler RP et al. 1979. Cytochrome C2 sequence
1914 variation among the recognised species of purple
1915 nonsulphur photosynthetic bacteria. *Nature*.
1916 278:659-660.

1917 Angly FE et al. 2014. CopyRighter: A rapid
1918 tool for improving the accuracy of microbial
1919 community profiles through lineage-specific gene
1920 copy number correction. *Microbiome*. 2:11.

1921 Apweiler R et al. 2004. UniProt: the Universal Pro-
1922 tein knowledgebase. *Nucleic Acids Res.* 32:D115-
1923 119.

1924 Arboleya S et al. 2018. Gene-trait matching
1925 across the *Bifidobacterium longum* pan-genome
1926 reveals considerable diversity in carbohydrate
1927 catabolism among human infant strains. *BMC*
1928 *Genomics*. 19:33.

1929 Armour CR, Nayfach S, Pollard KS, Sharpton TJ.
1930 2019. A Metagenomic Meta-analysis Reveals
1931 Functional Signatures of Health and Disease
1932 in the Human Gut Microbiome. *mSystems*.
1933 4:e00332-18.

1934 Ayling M, Clark MD, Leggett RM. 2019. New
1935 approaches for metagenome assembly with short
1936 reads. *Brief. Bioinform.* 00:1-11.

1937 Bäckhed F et al. 2015. Dynamics and stabilization of
1938 the human gut microbiome during the first year
1939 of life. *Cell Host Microbe*. 17:690-703.

1940 Barott KL et al. 2012. Microbial to reef scale inter-
1941 actions between the reef-building coral *Montastraea*
1942 *annularis* and benthic algae. *Proc. R. Soc. B*
1943 *Biol. Sci.* 279:1655-1664.

1944 Beiko RG, Harlow TJ, Ragan MA. 2005. Highways of
1945 gene sharing in prokaryotes. *Proc. Natl. Acad.*
1946 *Sci. USA.* 102:14332-14337.

1947 Bowers RM et al. 2017. Minimum information
1948 about a single amplified genome (MISAG) and
1949 a metagenome-assembled genome (MIMAG) of
1950 bacteria and archaea. *Nat. Biotechnol.* 35:725-
1951 731.

1952 Bowman JS, Ducklow HW. 2015. Microbial commu-
1953 nities can be described by metabolic structure: A
1954 general framework and application to a season-
1955 ally variable, depth-stratified microbial commu-
1956 nity from the coastal West Antarctic Peninsula.
1957 *PLOS One*. 10:e0135868.

1958 Bradley PH, Nayfach S, Pollard KS. 2018.
1959 Phylogeny-corrected identification of microbial
1960 gene families relevant to human gut colonization.
1961 *PLOS Comput. Biol.* 14:e1006242.

1962 Bradley PH, Pollard KS. 2020. Phylogenize: Correct-
1963 ing for phylogeny reveals genes associated with
1964 microbial distributions. *Bioinformatics*. 36:1289-
1965 1290.

Breiman L. 2001. Random Forests. *Mach. Learn.* 45:5-32. 1966

Brenner DJ. 1973. Deoxyribonucleic acid reassocia- 1967
tion in the taxonomy of enteric bacteria. *Int. J.* 1968
Syst. Bacteriol. 23:298-307. 1969

Buchfink B, Xie C, Huson DH. 2015. Fast and 1970
Sensitive Protein Alignment using DIAMOND. 1971
Nat. Methods. 12:59-60. 1972

Bukin YS et al. 2019. The effect of 16S rRNA region 1973
choice on bacterial community metabarcoding 1974
results. *Sci. Data*. 6:190007. 1975

Burke C, Steinberg P, Rusch DB, Kjelleberg S, 1976
Thomas T. 2011. Bacterial community assembly 1977
based on functional genes rather than species. 1978
Proc. Natl. Acad. Sci. USA. 108:14288-14293. 1979

Calgaro M, Romualdi C, Waldron L, Risso D, Vitulo 1980
N. 2020. Assessment of single cell RNA-seq 1981
statistical methods on microbiome data. *Genome* 1982
Biol. 191. 1983

Callahan BJ et al. 2016a. DADA2: High resolution 1984
sample inference from amplicon data. *Nat.* 1985
Methods. 13:581-583. 1986

Callahan BJ et al. 2019. High-throughput amplicon 1987
sequencing of the full-length 16S rRNA gene with 1988
single-nucleotide resolution. *Nucleic Acids Res.* 1989
47:e103. 1990

Callahan BJ, Sankaran K, Fukuyama JA, McMurdie 1991
PJ, Holmes SP. 2016b. Bioconductor workflow 1992
for microbiome data analysis: From raw reads to 1993
community analyses. *F1000 Res.* 5:1492. 1994

Caspi R et al. 2013. The MetaCyc database of 1995
metabolic pathways and enzymes and the BioCyc 1996
collection of pathway/genome databases. *Nucleic* 1997
Acids Res. 42:D459-D471. 1998

Castaner O et al. 2018. The gut microbiome 1999
profile in obesity: A systematic review. *Int. J.* 2000
Endocrinol. 2018:4095789. 2001

Chaffron S, Rehrauer H, Pernthaler J, Von Mering 2002
C. 2010. A global network of coexisting microbes 2003
from environmental and whole-genome sequence 2004
data. *Genome Res.* 20:947-959. 2005

Chen PE et al. 2010. Genomic characterization of 2006
the *Yersinia* genus. *Genome Biol.* 11:R1. 2007

Chen Z et al. 2019. Impact of Preservation Method 2008
and 16S rRNA Hypervariable Region on Gut 2009
Microbiota Profiling. *mSystems*. 4:e00271-18. 2010

Comeau AM, Douglas GM, Langille MGI. 2017. 2011
Microbiome Helper: a Custom and Streamlined 2012
Workflow for Microbiome Research. *mSystems*. 2013
2:e00127-16. 2014

Comeau AM, Lagunas MG, Scarcella K, Varela DE, 2015
Lovejoy C. 2019. Nitrate Consumers in Arctic 2016
Marine Eukaryotic Communities: Comparative 2017
Diversities of 18S rRNA, 18S rRNA Genes, and 2018
1916

2020 Nitrate Reductase Genes. *Appl. Environ. Microbiol.* 85:e00247-19. 2073

2021 2074

2022 Darling AE et al. 2014. PhyloSift: Phylogenetic 2075

2023 analysis of genomes and metagenomes. *PeerJ.* 2076

2024 2014:243. 2077

2025 Devereux R et al. 1990. Diversity and origin of 2078

2026 *Desulfovibrio* species: Phylogenetic definition of 2079

2027 a family. *J. Bacteriol.* 172:3609-3619. 2080

2028 Doolittle WF, Booth A. 2017. It's the song, not 2081

2029 the singer: an exploration of holobiosis and 2082

2030 evolutionary theory. *Biol. Philos.* 32:5-24. 2083

2031 Douglas GM et al. 2018. Multi-omics differentially 2084

2032 classify disease state and treatment outcome in 2085

2033 pediatric Crohn's disease. *Microbiome.* 6:13. 2086

2034 Douglas GM et al. 2020. PICRUST2 for prediction of 2087

2035 metagenome functions. *Nat. Biotechnol.* 38:685- 2088

2036 688. 2089

2037 Drouin G, Godin J-R, Pagé B. 2011. The genetics of 2090

2038 vitamin C loss in vertebrates. *Curr. Genomics.* 2091

2039 12:371-8. 2092

2040 Duvallet C, Gibbons SM, Gurry T, Irizarry RA, 2093

2041 Alm EJ. 2017. Meta-analysis of gut microbiome 2094

2042 studies identifies disease-specific and shared re- 2095

2043 sponses. *Nat. Commun.* 8:1784. 2096

2044 Egozcue JJ, Pawlowsky-Glahn V. 2011. Exploring 2097

2045 compositional data with the CoDa-dendrogram. 2098

2046 *Austrian J. Stat.* 40:103-113. 2099

2047 Eng A, Borenstein E. 2018. Taxa-function robustness 2100

2048 in microbial communities. *Microbiome.* 6:45. 2101

2049 Farrelly V, Rainey FA, Stackebrandt E. 1995. Effect 2102

2050 of genome size and rrn gene copy number on PCR 2103

2051 amplification of 16S rRNA genes from a mixture 2104

2052 of bacterial species. *Appl. Environ. Microbiol.* 2105

2053 61:2798-2801. 2106

2054 Faust K et al. 2012. Microbial co-occurrence 2107

2055 relationships in the Human Microbiome. *PLOS* 2108

2056 *Comput. Biol.* 8. 2109

2057 Fernandes AD et al. 2014. Unifying the analysis of 2110

2058 high-throughput sequencing datasets: character- 2111

2059 izing RNA-seq, 16S rRNA gene sequencing and 2112

2060 selective growth experiments by compositional 2113

2061 data analysis. *Microbiome.* 2:15. 2114

2062 Fernandes AD, Macklaim JM, Linn TG, Reid G, 2115

2063 Gloor GB. 2013. ANOVA-Like Differential Ex- 2116

2064 pression (ALDEx) Analysis for Mixed Population 2117

2065 RNA-Seq. *PLOS One.* 8:e67019. 2118

2066 Finn RD et al. 2014. Pfam: The protein families 2119

2067 database. *Nucleic Acids Res.* 42:D222-D230. 2120

2068 Fitch WM, Margoliash E. 1967. Construction of 2121

2069 phylogenetic trees. *Science.* 155:279-284. 2122

2070 Francke C, Siezen RJ, Teusink B. 2005. Reconstruct- 2123

2071 ing the metabolic network of a bacterium from 2124

2072 its genome. *Trends Microbiol.* 13:550-558. 2125

Franzosa EA et al. 2018. Species-level functional pro- 2073

2074 filing of metagenomes and metatranscriptomes. 2075

Nat. Methods. 15:962-968. 2076

Friedman J, Alm EJ. 2012. Inferring Correlation 2077

2078 Networks from Genomic Survey Data. *PLOS* 2079

Comput. Biol. 8:e1002687. 2080

Frioux C, Singh D, Korcsmaros T, Hildebrand F. 2081

2020. From bag-of-genes to bag-of-genomes: 2082

2083 metabolic modelling of communities in the era 2084

of metagenome-assembled genomes. *Comput. 2085*

Struct. Biotechnol. J. 18:1722-1734. 2086

Galperin MY, Kristensen DM, Makarova KS, Wolf 2087

YI, Koonin E V. 2019. Microbial genome anal- 2088

2089 ysis: The COG approach. *Brief. Bioinform.* 2090

20:1063-1070. 2091

Galperin MY, Makarova KS, Wolf YI, Koonin E V. 2092

2015. Expanded microbial genome coverage and 2093

2094 improved protein family annotation in the COG 2095

2096 database. *Nucleic Acids Res.* 43:D261-D269. 2097

Gevers D et al. 2014. The treatment-naive micro- 2098

2099 biome in new-onset Crohn's disease. *Cell Host 2100*

Microbe. 15:382-392. 2101

Gloor GB, Macklaim JM, Pawlowsky-Glahn V, 2102

Egozcue JJ. 2017. Microbiome datasets are 2103

2104 compositional: And this is not optional. *Front. 2105*

Microbiol. 8:2224. 2106

Gloor GB, Wu JR, Pawlowsky-Glahn V, Egozcue JJ. 2107

2016. It's all relative: analyzing microbiome data 2108

as compositions. *Ann. Epidemiol.* 26:322-329. 2109

Goodrich JK et al. 2014. Conducting a microbiome 2110

2111 study. *Cell.* 158:250-262. 2112

Graspeuntner S, Loeper N, Künzel S, Baines JF, 2113

Rupp J. 2018. Selection of validated hypervari- 2114

2115 able regions is crucial in 16S-based microbiota 2116

2117 studies of the female genital tract. *Sci. Rep.* 2118

8:9678. 2119

Greenblum S, Chiu HC, Levy R, Carr R, Boren- 2120

stein E. 2013. Towards a predictive systems- 2121

2122 level model of the human microbiome: Progress, 2123

2124 challenges, and opportunities. *Curr. Opin. 2125*

Biotechnol. 24:810-820.

Guo F, Zhang T. 2013. Biases during DNA ex- 2126

2127 traction of activated sludge samples revealed by 2128

2129 high throughput sequencing. *Appl. Microbiol. 2130*

Biotechnol. 97:4607-4616. 2131

Haft DH, Selengut JD, White O. 2003. The TIGR- 2132

2133 FAMS database of protein families. *Nucleic Acids 2134*

Res. 31:371-373. 2135

Hallam SJ, Girguis PR, Preston CM, Richardson 2136

PM, DeLong EF. 2003. Identification of Methyl 2137

2138 Coenzyme M Reductase A (*mcrA*) Genes Asso- 2139

2140 ciated with Methane-Oxidizing Archaea. *Appl. 2141*

Environ. Microbiol. 69:5483-5491. 2142

- 2126 Hao X, Chen T. 2012. OTU Analysis Using Metagenomic Shotgun Sequencing Data. *PLOS One*. 7:e49785. 2179
- 2127 2180
- 2128 2181
- 2129 Hauben L, Vauterin L, Moore ERB, Hoste B, Swings J. 1999. Genomic diversity of the genus *Stenotrophomonas*. *Int. J. Syst. Bacteriol.* 49:1749-1760. 2182
- 2130 2183
- 2131 2184
- 2132 2185
- 2133 Hauben L, Vauterin L, Swings J, Moore ERB. 1997. Comparison of 16S Ribosomal DNA Sequences of All *Xanthomonas* Species. *Int. J. Syst. Bacteriol.* 47:328-335. 2186
- 2134 2187
- 2135 2188
- 2136 2189
- 2137 Hawinkel S, Mattiello F, Bijmens L, Thas O. 2019. A broken promise: Microbiome differential abundance methods do not control the false discovery rate. *Brief. Bioinform.* 20:1-12. 2190
- 2138 2191
- 2139 2192
- 2140 2193
- 2141 Hill C. 2020. You have the microbiome you deserve. *Gut Microbiome*. 1:1-4. 2194
- 2142 2195
- 2143 Hillmann B et al. 2018. Evaluating the Information Content of Shallow Shotgun Metagenomics. *mSystems*. 3:e00069-18. 2196
- 2144 2197
- 2145 2198
- 2146 HMP-consortium. 2013. Structure, Function and Diversity of the Healthy Human Microbiome. *Nature*. 486:207-214. 2199
- 2147 2200
- 2148 2201
- 2149 Huson DH, Auch AF, Qi J, Schuster SC. 2007. MEGAN analysis of metagenomic data. *Genome Res.* 17:377-386. 2202
- 2150 2203
- 2151 2204
- 2152 Ibrahim A, Goebel BM, Liesack W, Griffiths M, Stackebrandt E. 1993. The phylogeny of the genus *Yersinia* based on 16S rDNA sequences. *FEMS Microbiol. Lett.* 114:173-177. 2205
- 2153 2206
- 2154 2207
- 2155 2208
- 2156 Inkpen AI et al. 2017. The Coupling of Taxonomy and Function in Microbiomes. *Biol. Philos.* 32:1225-1243. 2209
- 2157 2210
- 2158 2211
- 2159 Iwai S et al. 2016. Piphillin: Improved prediction of metagenomic content by direct inference from human microbiomes. *PLOS One*. 11:e0166104. 2212
- 2160 2213
- 2161 2214
- 2162 Jackson DA. 1997. Compositional data in community ecology: The paradigm or peril of proportions? *Ecology*. 78:929-940. 2215
- 2163 2216
- 2164 2217
- 2165 Janda JM, Abbott SL. 2007. 16S rRNA gene sequencing for bacterial identification in the diagnostic laboratory: Pluses, perils, and pitfalls. *J. Clin. Microbiol.* 45:2761-2764. 2218
- 2166 2219
- 2167 2220
- 2168 2221
- 2169 Jayaprakash TP, Schellenberg JJ, Hill JE. 2012. Resolution and characterization of distinct cpn60-based subgroups of *Gardnerella vaginalis* in the vaginal microbiota. *PLOS One*. 7:e43009. 2222
- 2170 2223
- 2171 2224
- 2172 2225
- 2173 Jensen LJ et al. 2008. eggNOG: Automated construction and annotation of orthologous groups of genes. *Nucleic Acids Res.* 36:D250-D254. 2226
- 2174 2227
- 2175 2228
- 2176 Johnson JS et al. 2019. Evaluation of 16S rRNA gene sequencing for species and strain-level microbiome analysis. *Nat. Commun.* 10:5029. 2229
- 2177 2230
- 2178 2231
- Jones MB et al. 2015. Library preparation methodology can influence genomic and functional predictions in human microbiome research. *Proc. Natl. Acad. Sci. USA*. 112:14024-14029. 2179
- 2180 2181
- 2182 2183
- 2183 2184
- 2184 2185
- 2185 2186
- 2186 2187
- 2187 2188
- 2188 2189
- 2189 2190
- 2190 2191
- 2191 2192
- 2192 2193
- 2193 2194
- 2194 2195
- 2195 2196
- 2196 2197
- 2197 2198
- 2198 2199
- 2199 2200
- 2200 2201
- 2201 2202
- 2202 2203
- 2203 2204
- 2204 2205
- 2205 2206
- 2206 2207
- 2207 2208
- 2208 2209
- 2209 2210
- 2210 2211
- 2211 2212
- 2212 2213
- 2213 2214
- 2214 2215
- 2215 2216
- 2216 2217
- 2217 2218
- 2218 2219
- 2219 2220
- 2220 2221
- 2221 2222
- 2222 2223
- 2223 2224
- 2224 2225
- 2225 2226
- 2226 2227
- 2227 2228
- 2228 2229
- 2229 2230
- 2230 2231

2232 Konstantinidis KT, Tiedje JM. 2005. Genomic 2286
2233 insights that advance the species definition for 2287
2234 prokaryotes. *Proc. Natl. Acad. Sci. USA.* 102:2567-2572. 2288
2235 2289
2236 Koonin E V, Galperin MY. 2003. Sequence - Evo- 2290
2237 lution - Function: Computational Approaches 2291
2238 in Comparative Genomics. Kluwer Academic: 2292
2239 Boston. 2293
2240 Kurtz ZD et al. 2015. Sparse and Compositionally 2294
2241 Robust Inference of Microbial Ecological Net- 2295
2242 works. *PLOS Comput. Biol.* 11:e1004226. 2296
2243 Langille MGI et al. 2013. Predictive functional 2297
2244 profiling of microbial communities using 16S 2298
2245 rRNA marker gene sequences. *Nat. Biotechnol.* 31:814-821. 2299
2246 2300
2247 Laserna-Mendieta EJ et al. 2018. Determinants of 2301
2248 reduced genetic capacity for butyrate synthesis 2302
2249 by the gut microbiome in Crohn's disease and 2303
2250 ulcerative colitis. *J. Crohn's Colitis.* 12:204-216. 2304
2251 Lau JT et al. 2016. Capturing the diversity of the 2305
2252 human gut microbiota through culture-enriched 2306
2253 molecular profiling. *Genome Med.* 8:72. 2307
2254 Ley RE, Peterson DA, Gordon JI. 2006. Ecological 2308
2255 and evolutionary forces shaping microbial diver- 2309
2256 sity in the human intestine. *Cell.* 124:837-848. 2310
2257 Li D, Liu CM, Luo R, Sadakane K, Lam TW. 2015. 2311
2258 MEGAHIT: An ultra-fast single-node solution 2312
2259 for large and complex metagenomics assembly 2313
2260 via succinct de Bruijn graph. *Bioinformatics.* 31:1674-1676. 2314
2261 2315
2262 Links MG, Dumonceaux TJ, Hemmingsen SM, Hill 2316
2263 JE. 2012. The Chaperonin-60 Universal Target 2317
2264 Is a Barcode for Bacteria That Enables De Novo 2318
2265 Assembly of Metagenomic Sequence Data. *PLOS* 2319
2266 *One.* 7:e49755. 2320
2267 Liu J, Yu Y, Cai Z, Bartlam M, Wang Y. 2015. 2321
2268 Comparison of ITS and 18S rDNA for estimating 2322
2269 fungal diversity using PCR-DGGE. *World J.* 2323
2270 *Microbiol. Biotechnol.* 31:1387-1395. 2324
2271 Lloyd-Price J et al. 2019. Multi-omics of the 2325
2272 gut microbial ecosystem in inflammatory bowel 2326
2273 diseases. *Nature.* 569:655-662. 2327
2274 Lloyd-Price J et al. 2017. Strains, functions and 2328
2275 dynamics in the expanded Human Microbiome 2329
2276 Project. *Nature.* 550:61-66. 2330
2277 Louca S et al. 2018a. Function and functional 2331
2278 redundancy in microbial systems. *Nat. Ecol.* 2332
2279 *Evol.* 2:936-943. 2333
2280 Louca S, Doebeli M. 2017. Taxonomic variability 2334
2281 and functional stability in microbial communities 2335
2282 infected by phages. *Environ. Microbiol.* 19:3863- 2336
2283 3878. 2337
2284 Louca S, Doebeli M, Parfrey LW. 2018b. Correcting 2338
2285 for 16S rRNA gene copy numbers in microbiome 2339
surveys remains an unsolved problem. *Micro-*
biome. 6:41.
Louca S, Parfrey LW, Doebeli M. 2016. Decoupling
function and taxonomy in the global ocean mi-
crobiome. *Science.* 353:1272-1277.
Love MI, Huber W, Anders S. 2014. Moderated esti-
mation of fold change and dispersion for RNA-seq
data with DESeq2. *Genome Biol.* 15:550.
Lozupone C, Knight R. 2005. UniFrac: a New
Phylogenetic Method for Comparing Microbial
Communities. *Appl. Environ. Microbiol.*
71:8228-8235.
Lozupone CA et al. 2008. The convergence of
carbohydrate active gene repertoires in human
gut microbes. *Proc. Natl. Acad. Sci. USA.*
105:15076-15081.
Lu J, Breitwieser FP, Thielen P, Salzberg SL.
2017. Bracken: Estimating species abundance in
metagenomics data. *PeerJ Comput. Sci.* 3:e104.
Makarova KS, Wolf YI, Koonin E V. 2015. Archaeal
clusters of orthologous genes (arCOGs): An
update and application for analysis of shared fea-
tures between thermococcales, methanococcales,
and methanobacteriales. *Life.* 5:818-840.
Mandal S et al. 2015. Analysis of composition
of microbiomes: a novel method for studying
microbial composition. *Microb. Ecol. Heal. Dis.*
26:27663.
Mandel M. 1966. Deoxyribonucleic Acid Base Com-
position in the Genus *Pseudomonas*. *J. Gen.*
Microbiol. 43:273-292.
Manor O, Borenstein E. 2017a. Revised computa-
tional metagenomic processing uncovers hidden
and biologically meaningful functional variation
in the human microbiome. *Microbiome.* 5:19.
Manor O, Borenstein E. 2017b. Systematic Charac-
terization and Analysis of the Taxonomic Drivers
of Functional Shifts in the Human Microbiome.
Cell Host Microbe. 21:254-267.
Martin BD, Witten D, Willis AD. 2020. Modeling
microbial abundances and dysbiosis with beta-
binomial regression. *Ann. Appl. Stat.* 14:94-
115.
Martiny AC. 2019. High proportions of bacteria
are culturable across major biomes. *ISME J.*
13:2125-2128.
Martiny AC, Treseder K, Pusch G. 2013. Phyloge-
netic conservatism of functional traits in microor-
ganisms. *ISME J.* 7:830-838.
Maukonen J, Simões C, Saarela M. 2012. The cur-
rently used commercial DNA-extraction methods
give different results of clostridial and actinobac-
terial populations derived from human fecal sam-
ples. *FEMS Microbiol. Ecol.* 79:697-708.

2340 McDonald D et al. 2019. redbiom: a Rapid Sample
2341 Discovery and Feature Characterization System.
2342 mSystems. 4:e00215-19. 2394

2343 McIntyre ABR et al. 2017. Comprehensive bench-
2344 marking and ensemble approaches for metage-
2345 nomic classifiers. Genome Biol. 18:182. 2395

2346 McMahon K. 2015. "Metagenomics 2.0". Environ.
2347 Microbiol. Rep. 7:38-39. 2396

2348 Meyer F et al. 2008. The metagenomics RAST server
2349 - A public resource for the automatic phyloge-
2350 netic and functional analysis of metagenomes.
2351 BMC Bioinformatics. 9:386. 2397

2352 Meyer F, Overbeek R, Rodriguez A. 2009. FIGfams:
2353 Yet another set of protein families. Nucleic Acids
2354 Res. 37:6643-6654. 2398

2355 Miossec MJ et al. 2020. Evaluation of computational
2356 methods for human microbiome analysis using
2357 simulated data. PeerJ. 8:e9688. 2399

2358 Morgan XC et al. 2012. Dysfunction of the intestinal
2359 microbiome in inflammatory bowel disease and
2360 treatment. Genome Biol. 13:R79. 2400

2361 Moriya Y, Itoh M, Okuda S, Yoshizawa AC, Kanehisa
2362 M. 2007. KAAS: An automatic genome annota-
2363 tion and pathway reconstruction server. Nucleic
2364 Acids Res. 35:W182-W185. 2401

2365 Morris JJ, Lenski RE, Zinser ER. 2012. The Black
2366 Queen Hypothesis: Evolution of Dependencies
2367 through Adaptive Gene Loss. MBio. 3:e00036-
2368 12. 2402

2369 Morton JT et al. 2017. Balance Trees Reveal Micro-
2370 bial Niche Differentiation. mSystems. 2:e00162-
2371 16. 2403

2372 Morton JT et al. 2019. Establishing microbial com-
2373 position measurement standards with reference
2374 frames. Nat. Commun. 10:2719. 2404

2375 Mosca A, Leclerc M, Hugot JP. 2016. Gut microbiota
2376 diversity and human diseases: Should we rein-
2377 troduce key predators in our ecosystem? Front.
2378 Microbiol. 7:455. 2405

2379 Moya A, Ferrer M. 2016. Functional Redundancy-
2380 Induced Stability of Gut Microbiota Subjected
2381 to Disturbance. Trends Microbiol. 24:402-413. 2406

2382 Muegge BD et al. 2011. Diet Drives Conver-
2383 gence in Gut Microbiome Functions Across Mam-
2384 malian Phylogeny and Within Humans. Science.
2385 332:970-974. 2407

2386 Muscogiuri G et al. 2019. Gut microbiota: a new
2387 path to treat obesity. Int. J. Obes. Suppl. 9:10-
2388 19. 2408

2389 Mysara M et al. 2017. Reconciliation between oper-
2390 ational taxonomic units and species boundaries.
2391 FEMS Microbiol. Ecol. 93:fix029. 2409

2392 Naeem S, Kawabata Z, Loreau M. 1998. Transcend-
2393 ing boundaries in biodiversity research. Trends
2394 Ecol. Evol. 13:134-135. 2395

2396 Narayan NR et al. 2020. Piphillin predicts metage-
2397 nomic composition and dynamics from DADA2-
2398 corrected 16S rDNA sequences. BMC Genomics.
2399 21:56. 2400

2400 Nearing JT, Douglas GM, Comeau AM, Langille
2401 MGI. 2018. Denoising the Denoisers: An in-
2402 dependent evaluation of microbiome sequence
2403 error- correction approaches. PeerJ. 2018:e5364. 2404

2404 Nejman D et al. 2020. The human tumor microbiome
2405 is composed of tumor type-specific intra-cellular
2406 bacteria. Science. 980:973-980. 2407

2407 NIH. 2019. A review of 10 years of human mi-
2408 crobiome research activities at the US National
2409 Institutes of Health, Fiscal Years 2007-2016.
2410 Microbiome. 7:31. 2408

2411 Ning J, Beiko RG. 2015. Phylogenetic approaches to
2412 microbial community classification. Microbiome.
2413 3:47. 2409

2414 Nurk S, Meleshko D, Korobeynikov A, Pevzner PA.
2415 2017. metaSPAdes: a new versatile metagenomic
2416 assembler. Genome Res. 27:824-834. 2410

2417 Oberhardt MA, Pucha?ka J, Fryer KE, Martins Dos
2418 Santos VAP, Papin JA. 2008. Genome-scale
2419 metabolic network analysis of the opportunistic
2420 pathogen *Pseudomonas aeruginosa* PAO1. J.
2421 Bacteriol. 190:2790-2803. 2411

2422 Oh J et al. 2014. Biogeography and individuality
2423 shape function in the human skin metagenome.
2424 Nature. 514:59-64. 2412

2425 Olson ND et al. 2019. Metagenomic assembly
2426 through the lens of validation: Recent advances
2427 in assessing and improving the quality of genomes
2428 assembled from metagenomes. Brief. Bioinform.
2429 20:1140-1150. 2413

2430 Omelchenko M V., Galperin MY, Wolf YI, Koonin E
2431 V. 2010. Non-homologous isofunctional enzymes:
2432 A systematic analysis of alternative solutions in
2433 enzyme evolution. Biol. Direct. 5:31. 2414

2434 Parks DH, Imelfort M, Skennerton CT, Hugenholtz P,
2435 Tyson GW. 2015. CheckM: assessing the quality
2436 of microbial genomes recovered from isolates,
2437 single cells, and metagenomes. Genome Res.
2438 25:1043-55. 2415

2439 Pasolli E, Truong T, Malik F, Waldron L, Segata
2440 N. 2016. Machine learning meta-analysis of
2441 large metagenomic datasets: tools and biological
2442 insights. PLOS Comput. Biol. 12:e1004977. 2416

2443 Paulson JN, Colin Stine O, Bravo HC, Pop M.
2444 2013. Differential abundance analysis for mi-
2445 crobial marker-gene surveys. Nat. Methods.
2446 10:1200-1202. 2417

2447 Pollock J, Glendinning L, Wisedchanwet T, Watson
2448 M. 2018. The Madness of Microbiome: Attempt-

ing To Find Consensus “Best Practice” for 16S
Microbiome Studies. *Appl. Environ. Microbiol.*
84:e02627-17.

2451 Popa O, Dagan T. 2011. Trends and barriers to
2452 lateral gene transfer in prokaryotes. *Curr. Opin.*
2453 *Microbiol.* 14:615-623.

2454 Prakash T, Taylor TD. 2012. Functional assignment
2455 of metagenomic data: challenges and applica-
2456 tions. *Brief. Bioinform.* 13:711-727.

2457 Prodan A et al. 2020. Comparing bioinformatic
2458 pipelines for microbial 16S rRNA amplicon se-
2459 quencing. *PLOS One.* 15:e0227434.

2460 Punta M et al. 2012. The Pfam protein families
2461 database. *Nucleic Acids Res.* 40:D290-D301.

2462 Quast C et al. 2013. The SILVA ribosomal RNA gene
2463 database project: Improved data processing and
2464 web-based tools. *Nucleic Acids Res.* 41:590-596.

2465 Rahman SF, Olm MR, Morowitz MJ, Banfield JF.
2466 2018. Machine Learning Leveraging Genomes
2467 from Metagenomes Identifies Influential Anti-
2468 biotic Resistance Genes in the Infant Gut Micro-
2469 biome. *mSystems.* 3:e00123-17.

2470 Rasko DA et al. 2008. The pangenome structure of
2471 *Escherichia coli*: Comparative genomic analysis
2472 of *E. coli* commensal and pathogenic isolates. *J.*
2473 *Bacteriol.* 190:6881-6893.

2474 Riley M. 1993. Functions of the gene products of
2475 *Escherichia coli*. *Microbiol. Rev.* 57:862-952.

2476 Salonen A et al. 2010. Comparative analysis of
2477 fecal DNA extraction methods with phylogenetic
2478 microarray: Effective recovery of bacterial and
2479 archaeal DNA using mechanical cell lysis. *J.*
2480 *Microbiol. Methods.* 81:127-134.

2481 Saroj DB, Dengeti SN, Aher S, Gupta AK. 2015. ITS
2482 as an environmental DNA barcode for fungi: an
2483 *in silico* approach reveals potential PCR biases.
2484 *World J. Microbiol. Biotechnol.* 31:189.

2485 Schloss PD. 2020. Reintroducing mothur: 10 Years
2486 Later. *Appl. Environ. Microbiol.* 86:e02343-19.

2487 Schoch CL et al. 2012. Nuclear ribosomal internal
2488 transcribed spacer (ITS) region as a universal
2489 DNA barcode marker for Fungi. *Proc. Natl.*
2490 *Acad. Sci. USA.* 109:6241-6246.

2491 Schwager E, Mallick H, Ventz S, Huttenhower C.
2492 2017. A Bayesian method for detecting pairwise
2493 associations in compositional data. *PLOS Com-
2494 put. Biol.* 13:e1005852.

2495 Segata N et al. 2011. Metagenomic biomarker dis-
2496 covery and explanation. *Genome Biol.* 12:R60.

2497 Shade A. 2017. Diversity is the question, not the
2498 answer. *ISME J.* 11:1-6.

2499 Shapley LS. 1953. A value for n-person games. In:
2500 *Contributions to the Theory of Games, 2.* Kuhn,
HW & Tucker, W, editors. Princeton University
Press: Princeton, NJ pp. 307-317.

Silverman JD, Washburne AD, Mukherjee S, David
LA. 2017. A phylogenetic transform enhances
analysis of compositional microbiota data. *Elife.*
6:e21887.

Sperling JL et al. 2017. Comparison of bacterial 16S
rRNA variable regions for microbiome surveys of
ticks. *Ticks Tick. Borne. Dis.* 8:453-461.

Sprockett D et al. 2019. Treatment-specific com-
position of the gut microbiota is associated with
disease remission in a pediatric Crohn’s disease
cohort. *Inflamm. Bowel Dis.* 25:1927-1938.

Stackebrandt E, Goebel BM. 1994. Taxonomic note:
A place for DNA-DNA reassociation and 16S
rRNA sequence analysis in the present species
definition in bacteriology. *Int. J. Syst. Bacteriol.*
44:846-849.

Staley J, Konopka A. 1985. Measurement of In Situ
Activities of Nonphotosynthetic Microorganisms
in Aquatic and Terrestrial Habitats. *Annu. Rev.*
Microbiol. 39:321-346.

Stein JL, Marsh TL, Wu KY, Shizuya H, DeLong EF.
1996. Characterization of uncultivated prokary-
otes: Isolation and analysis of a 40-kilobase-
pair genome fragment from a planktonic marine
archaeon. *J. Bacteriol.* 178:591-599.

Steinegger M, Söding J. 2018. Clustering huge pro-
tein sequence sets in linear time. *Nat. Commun.*
9:2542.

Steinegger M, Söding J. 2017. MMseqs2 enables sen-
sitive protein sequence searching for the analysis
of massive data sets. *Nat. Biotechnol.* 35:1026-
1028.

Suardana IW. 2014. Analysis of Nucleotide Se-
quences of the 16S rRNA Gene of Novel *Es-
cherichia coli* Strains Isolated from Feces of
Human and Bali Cattle. *J. Nucleic Acids.*
2014:475754.

Sun DL, Jiang X, Wu QL, Zhou NY. 2013. Intra-
genomic heterogeneity of 16S rRNA genes causes
overestimation of prokaryotic diversity. *Appl.*
Environ. Microbiol. 79:5962-5969.

Sun S, Jones RB, Fodor AA. 2020. Inference-based
accuracy of metagenome prediction tools varies
across sample types and functional categories.
Microbiome. 8:46.

Sze MA, Schloss PD. 2016. Looking for a signal in
the noise: Revisiting obesity and the microbiome.
MBio. 7:e01018-16.

Tatusov RL, Galperin MY, Natale DA, Koonin EV.
2000. The COG database: a tool for genome-
scale analysis of protein functions and evolution.
Nucleic Acids Res. 28:33-36.

- 2555 Tedjo DI et al. 2016. The fecal microbiota as a
2556 biomarker for disease activity in Crohn’s disease.
2557 Sci. Rep. 6:35216.
- 2558 Tessler M et al. 2017. Large-scale differences
2559 in microbial biodiversity discovery between 16S
2560 amplicon and shotgun sequencing. Sci. Rep.
2561 7:6589.
- 2562 Tettelin H et al. 2005. Genome analysis of multiple
2563 pathogenic isolates of *Streptococcus agalactiae*:
2564 Implications for the microbial “pan-genome”.
2565 Proc. Natl. Acad. Sci. USA. 102:3950-13955.
- 2566 Thompson LR et al. 2017. A communal catalogue
2567 reveals Earth’s multiscale microbial diversity.
2568 Nature. 551:457-463.
- 2569 Thorsen J et al. 2016. Large-scale benchmarking
2570 reveals false discoveries and count transformation
2571 sensitivity in 16S rRNA gene amplicon data
2572 analysis methods used in microbiome studies.
2573 Microbiome. 4:62.
- 2574 Treem W, Ahsan N, M S, Hyams J. 1994. Fecal
2575 Short-Chain Fatty Acids in Children with Inflam-
2576 matory Bowel Disease. J. Pediatr. Gastroen-
2577 terol. Nutr. 18:159-164.
- 2578 Truong DT et al. 2015. MetaPhlan2 for enhanced
2579 metagenomic taxonomic profiling. Nat. Meth-
2580 ods. 12:902-903.
- 2581 Turnbaugh PJ et al. 2009. A core gut microbiome in
2582 obese and lean twins. Nature. 457:480-484.
- 2583 Turnbaugh PJ, Bäckhed F, Fulton L, Gordon JI.
2584 2008. Diet-Induced Obesity Is Linked to Marked
2585 but Reversible Alterations in the Mouse Distal
2586 Gut Microbiome. Cell Host Microbe. 3:213-223.
- 2587 Vatanen T et al. 2018. Genomic variation and strain-
2588 specific functional adaptation in the human gut
2589 microbiome during early life. Nat. Microbiol.
2590 4:470-479.
- 2591 Venegas DP et al. 2019. Short chain fatty acids
2592 (SCFAs) mediated gut epithelial and immune
2593 regulation and its relevance for inflammatory
2594 bowel diseases. Front. Immunol. 10:277.
- 2595 Venter JC et al. 2004. Environmental Genome
2596 Shotgun Sequencing of the Sargasso Sea. Science.
2597 304:66-74.
- 2598 Verster AJ, Borenstein E. 2018. Competitive lottery-
2599 based assembly of selected clades in the human
2600 gut microbiome. Microbiome. 6:186.
- 2601 (Větrovský T, Baldrian P. 2013. The Variability
2602 of the 16S rRNA Gene in Bacterial Genomes
2603 and Its Consequences for Bacterial Community
2604 Analyses. PLOS One. 8:e57923.
- 2605 Vincent C et al. 2013. Reductions in intestinal
2606 Clostridiales precede the development of nosoco-
2607 mial *Clostridium difficile* infection. Microbiome.
2608 1:18.
- 2609 Wang Y, Qian P, Ya S. 2013. Conserved
2610 Regions in 16S Ribosome RNA Sequences
2611 and Primer Design for Studies of Envi-
2612 ronmental Microbes. Encycl. Metage-
2613 nomics. [https://doi.org/10.1007/978-1-4614-
2614 6418-1-772-1](https://doi.org/10.1007/978-1-4614-6418-1-772-1).
- 2615 Watson EJ, Giles J, Scherer BL, Blatchford P. 2019.
2616 Human faecal collection methods demonstrate
2617 a bias in microbiome composition by cell wall
2618 structure. Sci. Rep. 9:16831.
- 2619 Weiss S et al. 2017. Normalization and microbial
2620 differential abundance strategies depend upon
2621 data characteristics. Microbiome. 5:27.
- 2622 Wemheuer F et al. 2020. Tax4Fun2: prediction of
2623 habitat-specific functional profiles and functional
2624 redundancy based on 16S rRNA gene sequences.
2625 Environ. Microbiome. 15:11.
- 2626 Wheeler NE, Barquist L, Kingsley RA, Gardner PP.
2627 2016. A profile-based method for identifying
2628 functional divergence of orthologous genes in
2629 bacterial genomes. Bioinformatics. 32:3566-
2630 3574.
- 2631 Willis C, Desai D, Laroche J. 2019. Influence of
2632 16S rRNA variable region on perceived diversity
2633 of marine microbial communities of the North-
2634 ern North Atlantic. FEMS Microbiol. Lett.
2635 366:fnz152.
- 2636 Wilson GA et al. 2005. Orphans as taxonomically
2637 restricted and ecologically important genes. Mi-
2638 crobiology. 151:2499-2501.
- 2639 Wirbel J et al. 2019. Meta-analysis of fecal
2640 metagenomes reveals global microbial signatures
2641 that are specific for colorectal cancer. Nat. Med.
2642 25:679-689.
- 2643 Woese CR. 1987. Bacterial evolution. Microbiol.
2644 Rev. 51:221-271.
- 2645 Woese CR et al. 1980. Secondary structure model
2646 for bacterial 16S ribosomal RNA: Phylogenetic,
2647 enzymatic and chemical evidence. Nucleic Acids
2648 Res. 8:2275-2294.
- 2649 Woese CR, Fox GE. 1977. Phylogenetic structure of
2650 the prokaryotic domain: The primary kingdoms.
2651 Proc. Natl. Acad. Sci. USA. 74:5088-5090.
- 2652 Wood DE, Lu J, Langmead B. 2019. Improved
2653 metagenomic analysis with Kraken 2. Genome
2654 Biol. 20:257.
- 2655 Wright EK et al. 2015. Recent advances in character-
2656 izing the gastrointestinal microbiome in Crohn’s
2657 disease: a systematic review. Inflamm Bowel Dis.
2658 21:1219-1228.
- 2659 Wu D, Jospin G, Eisen JA. 2013. Systematic Identifi-
2660 cation of Gene Families for Use as ‘Markers’ for
2661 Phylogenetic and Phylogeny-Driven Ecological
2662 Studies of Bacteria and Archaea and Their Major

2663 Subgroups. PLOS One. 8:e77033.
2664 Ye SH, Siddle KJ, Park DJ, Sabeti PC. 2019.
2665 Benchmarking Metagenomics Tools for Taxo-
2666 nomic Classification. Cell. 178:779-794.
2667 Ye Y, Doak TG. 2011. A Parsimony Approach to
2668 Biological Pathway Reconstruction/Inference for
2669 Metagenomes. PLOS Comput. Biol. 5:e1000465.
2670 Zaneveld JR, Lozupone C, Gordon JI, Knight R.
2671 2010. Ribosomal RNA diversity predicts genome
2672 diversity in gut bacteria and their relatives.
2673 Nucleic Acids Res. 38:3869-3879.
2674 Zaneveld JRR, Thurber RLV. 2014. Hidden state
2675 prediction: A modification of classic ancestral
2676 state reconstruction algorithms helps unravel
2677 complex symbioses. Front. Microbiol. 5:431.
2678 Zhou J et al. 2015. High-Throughput Metagenomic
2679 Technologies for Complex Microbial Community
2680 Analysis: Open and Closed Formats. MBio.
2681 6:e02288-14.
2682 Zhou W et al. 2019. Longitudinal multi-omics of
2683 host-microbe dynamics in prediabetes. Nature.
2684 569:663-671.
2685 Zhou YH, Gallins P. 2019. A review and tutorial of
2686 machine learning methods for microbiome host
2687 trait prediction. Front. Genet. 10:579.
2688 Zuckerkandl E, Pauling L. 1965. Molecules as docu-
2689 ments of history. J. Theor. Biol. 8:357-366.

Reviewer 1

Dear Gavin Douglas and Morgan Langille,

In this review manuscript, you propose to deliver a detailed introduction to microbiome DNA sequence data types and analysis methods. You present marker-gene and shotgun DNA sequencing data types, discuss microbiome data characteristics and underscore the associated caveats. Then you present « the many-faceted concept of microbial functions ». You follow by a discussion on the problematic of functional annotation inferred from marker-gene data and you review the last development on the integration of taxonomy and function. Finally, you discuss reproducibility in microbiome research and provide an outlook with some personal experience.

The main strength of your manuscript is that as a reader I learned something because you deliver an interesting review and discussion on the integration of taxonomic and functional microbiome data, backed up by first-hand and authoritative experience. However, the main weakness is that your message is diluted by a lengthy and unclear explanation of some concepts that are not always directly linked to your main discussion point.

Therefore, I recommend a major revision of your manuscript.

Sincerely,

Nicolas Pollet

Major comments:

What is the audience ? The title says « A primer and discussion on DNA-based microbiome data and related bioinformatics analyses ». Since one aim is to deliver a primer, the reader is expected to be a non-expert, and therefore the discussion that follows is also expected to reach a non-expert in the field. Is this the case ? I don't think so. In fact, I am unsure of the efficiency to pursue the goal of fulfilling the role of a primer AND a discussion on microbiome data for a reader completely new to the field and for the complicated topics presented here. Since the reader could be misled by the title, I think you need to change it to better represent the content of the text.

Is the communication clear enough for a newcomer ? I think that you have to work on making your text more concise and more homogenous in terms of the depth of explanations. More and better iconography would help in this regard. The iconography should follow the main organization of the text : here you have six sections and only two figures. Figure 1 illustrates many aspects of the section on shotgun metagenomics, and figure 2 is an illustration on the integration of taxonomy and function. Figure 1d does not follow the text flow and I find this a bit strange.

What is the review message ? In my opinion, the discussion on the integration of taxonomic and functional data is the main message. I advise you to strengthen this aspect by dropping some sections (see below).

How to make the message clearer ? If you decide to follow the path of considering the integration of taxonomic and functional data as the main message to deliver, then the text could be reorganized to make this message stronger and clearer. I wonder if the sometime high level of details provided regarding marker-gene sequencing, shotgun metagenomics and the characteristics of microbiome count data is really helping the reader. The text would benefit from being way more concise and more equilibrated among sections. In my opinion, you should seriously consider to skip the “primer” sections on marker-gene sequencing, metagenomic sequencing, characteristics of microbiome count data and microbial functions.

I found that the discussion is the best part of the text, maybe because I am not a complete newcomer to the field. Your personal account is worthy, and maybe you could make it more precise (e.g. parameter choice from local to global using which tool?). The last two sections are the most informative parts and in this regard.

Accuracy : The terminology about microbiome is sound and corresponds to what has been previously discussed in the literature (Marchesi & Ravel, 2015). I found that the terminology used in the section microbial function is not always clear and does not simplify the presentation of the associated concepts (Karp, 2000)(Thomas, Mi & Lewis, 2007)(Kotera et al., 2014).

Level of referencing : There are specific experimental approaches such as epicPCR that have been developed to tackle the integration of taxonomy and function; and this needs to be pointed out (Spencer et al., 2016). I think you should take a particular attention to be more homogeneous in the way you select the cited references.

Minor comments

Since the review aims to deliver a detailed introduction, I suggest to expand a bit the terminology and definition that you provide rapidly for the term microbiome (one sentence on line 31-33), and possibly include a text-box with definitions. Maybe the ecological suffix -biome that refers to biotic and abiotic factors characterizing a given microbiome environment would broaden the scope.

I fully understand that the topic is DNA-based sequencing for microbiome studies, but a pointer to RNA-based and protein-based sequencing would be a plus in the background, especially in the paragraph 45-67. In that same paragraph on culturing microbes, and given the theme of the integration between taxonomy and function, one possible additional point could be to discuss the discrimination of live, dormant and dead microorganisms (e.g. (Thomas, Mi & Lewis, 2007) (Jones & Lennon, 2010)(Carini et al., 2016)(Blazewicz et al., 2013).

In the background section presenting diversity analysis, I would like to underscore the work of Amy Willis and colleagues on modelling abundances as in my opinion it is an important advance in the analysis of diversity (Willis, 2019)(Willis & Martin). The purpose of this paragraph in the context of the review as a whole is unclear as it stands.

I do not agree with the assertion that the dichotomy between phylogenetic and functional profiling of microbiomes is « entirely related to methodological challenges » (line 123). We know that the genome of prokaryotic species varies in gene content because of horizontal gene transfer, gene duplication and other mechanisms (Puigbò et al., 2014). It has been shown through pangenome analysis that strain variation can be associated with different metabolic potential (Goyal, 2018) (Maistrenko et al., 2020). Therefore, it seems to me that the dichotomy between phylogenetic and functional profiling of microbiomes is one of their intrinsic characteristics. Indeed, you develop these points line 1131-1172.

Marker-gene sequencing

I advise to simplify the marker gene sequencing section if you want to keep it. While the paragraph from 149-202 are detailed and very informative, I am afraid that they depart from the global « granularity » of explanation and historical context provided on other aspects throughout the manuscript. This lengthen this section on marker genes comparatively to the other aspects developed in this review. And even if there are a lot of things to tell about 16S rRNA gene sequencing, many have already been told elsewhere in the literature.

While I typically enjoy reading historical perspectives, I found that these are exaggeratedly long and placed in the manuscript in a non-logical manner.

You copiously present 16S rRNA gene sequencing and this helps the reader for understanding the aspects on the integration of taxonomic and functional data. But you also consider other marker genes (and this is fine) and 18s rRNA gene sequencing for microeukaryote and fungi taxonomic profiling, but in a more concise manner. Yet the integration of taxonomic data obtained using such markers with shotgun sequencing data is not presented at all, and thus the reader does not benefit from this otherwise interesting piece of knowledge.

The sentence line 211 would benefit from some simplification such as :

« This is because if there are non-random substitutions within a single domain but random substitutions in the majority of other domains, there would likely be little effect on estimates of gene divergence. »

I do not understand the reason for presenting redbiom at this point line 250 ?

To further document your point on the limitations due to the use of short 16S amplicons (line 260-274), you could possibly cite the recent work of other groups such as (Abellan-Schneyder et al., 2021).

The point dealing with the use of classical bacteria 16S primer-pairs do characterize Archaea could be expanded as it is often a neglected limitation in taxonomic surveys (Raymann et al., 2017; Bahram et al., 2019).

The reference Fox et al 1992 is missing at line 235. I think it would be fair to reference deblur and UNOISE3 like it has been made for DADA2 software (line 336).

Very Minor : italicize latin names (e.g Haloarcula line 382)

Shotgun metagenomics sequencing

Line 409 : including DNA viruses

The impact of biomass and genome size as a limitation to MGS approach could be invoked (line 431). Also as a caveat emptor, the impact of host DNA and possible heterologous sequences on MGS data could be mentioned, (I wrote this sentence before reading your discussion !) and this would be a reflection of the discussion.

In the MGS data analysis section devoted to the generation of taxonomic profile (line 477-522) , I would like to point out the targeted assembly of rRNA sequences from shotgun data embodied in Emirge (Miller et al., 2011), phyloFlash (Gruber-Vodicka, Seah & Pruesse, 2020) and MATAM (Pericard et al., 2018).

I was surprised that the authors do not mention Kaiju as a read-based tool for taxonomic profiling (Menzel, Ng & Krogh, 2016).

On the impact of databases for k-mer based analysis (Nasko et al., 2018).

Line 560 : the citation of only these two assemblers is somehow partial, you could point to a review on metagenome assembly for the sake of comprehensiveness for the reader. Similarly the description of binning tools is very light in comparison to other aspects developed earlier. Here you could point to recent review papers on the subject.

Line 584 : maybe use « taxonomic profiling » instead of « profiling »

Line 586 : I guess that the authors are referring to transcriptome studies, the term RNA sequencing is maybe not so precise in this context.

Characteristics of microbiome count data :

Maybe at some point the word abundance table could be used.

Line 618-637 : Maybe a figure would be a better communication vector.

The impact of sequencing reads processing on the analysis of abundance tables is somehow skipped : there are different practices such as removing singletons, filtering on prevalence etc . This could be somehow mentioned as they impact downstream analysis.

Microbial functions

This section is quite lengthy in comparison to others and since it covers topics that are not specific to microbiome studies, I wonder if it hits the sweet spot.

Line 737 : « ... focused on gene families, which are gene clusters. » It is not very clear what you are referring to in terms of gene cluster at this point.

Line 781 : I do not know what is a UniRef function.

What is described in this paragraph entitled microbial function is in fact a primer on protein databases and ontologies. I find therefore that the title is a bit misleading, maybe « Protein databases and ontologies for microbial genome functional annotation ».

Line 976 : this method focuses pathway reconstruction ... please correct the sentence.

Line 1032 : philosophical perspective : really ?

Line 1060 : The whole presentation of this paragraph is somehow paradoxal : maybe the text could be more explicit on ontology and semantics in order to guide the analysis of « functional data » at a given level of an ontology (protein space, biochemical activity, pathway, evolutionary conservation).

Metagenome prediction methods

Line 1090-1097 : some references would be welcome here.

Lines 1101-1110 -1130: This historical account is perfect, but I wonder if the level of details provided is really needed to make the point that 16S diversity is not a perfect proxy of whole genome similarity.

Current state of the integration of taxonomic and functional data types

I enjoyed reading this section.

Line 1313: “in some cases can be directly linked” Please be more precise and provide an example or a reference.

Why the burrito software is not mentioned is unclear to me ?

Outlook

In my opinion, the paragraph 1726-1761 would benefit from citing additional recent references such as the MBQC study and a few others: (Sinha et al., 2017; Davis et al., 2018; McLaren, Willis & Callahan, 2019; Greathouse, Sinha & Vogtmann, 2019).

References

I suggest to use a style for references that includes a DOI.

Abellan-Schneyder I, Matchado MS, Reitmeier S, Sommer A, Sewald Z, Baumbach J, List M, Neuhaus K. 2021. Primer, Pipelines, Parameters: Issues in 16S rRNA Gene Sequencing. *mSphere* 6. DOI: 10.1128/mSphere.01202-20.

Bahram M, Anslan S, Hildebrand F, Bork P, Tedersoo L. 2019. Newly designed 16S rRNA metabarcoding primers amplify diverse and novel archaeal taxa from the environment. *Environmental Microbiology Reports* 11:487–494. DOI: 10.1111/1758-2229.12684.

Blazewicz SJ, Barnard RL, Daly RA, Firestone MK. 2013. Evaluating rRNA as an indicator of microbial activity in environmental communities: limitations and uses. *The ISME journal* 7:2061–2068. DOI: 10.1038/ismej.2013.102.

Carini P, Marsden PJ, Leff JW, Morgan EE, Strickland MS, Fierer N. 2016. Relic DNA is abundant in soil and obscures estimates of soil microbial diversity. *Nature Microbiology* 2:1–6. DOI: 10.1038/nmicrobiol.2016.242.

Davis NM, Proctor DM, Holmes SP, Relman DA, Callahan BJ. 2018. Simple statistical identification and removal of contaminant sequences in marker-gene and metagenomics data. *Microbiome* 6:226. DOI: 10.1186/s40168-018-0605-2.

Emerson JB, Adams RI, Román CMB, Brooks B, Coil DA, Dahlhausen K, Ganz HH, Hartmann EM, Hsu T, Justice NB, Paulino-Lima IG, Luongo JC, Lympelopoulou DS, Gomez-Silvan C, Rothschild-Mancinelli B, Balk M, Huttenhower C, Nocker A, Vaishampayan P, Rothschild LJ. 2017. Schrödinger's microbes: Tools for distinguishing the living from the dead in microbial ecosystems. *Microbiome* 5:86. DOI: 10.1186/s40168-017-0285-3.

Goyal A. 2018. Metabolic adaptations underlying genome flexibility in prokaryotes. *PLOS Genetics* 14:e1007763. DOI: 10.1371/journal.pgen.1007763.

Greathouse KL, Sinha R, Vogtmann E. 2019. DNA extraction for human microbiome studies: the issue of standardization. *Genome Biology* 20:212. DOI: 10.1186/s13059-019-1843-8.

Gruber-Vodicka HR, Seah BKB, Pruesse E. 2020. phyloFlash: Rapid Small-Subunit rRNA Profiling and Targeted Assembly from Metagenomes. *mSystems* 5. DOI: 10.1128/mSystems.00920-20.

Jones SE, Lennon JT. 2010. Dormancy contributes to the maintenance of microbial diversity. *Proceedings of the National Academy of Sciences* 107:5881–5886. DOI: 10.1073/pnas.0912765107.

Karp PD. 2000. An ontology for biological function based on molecular interactions. *Bioinformatics (Oxford, England)* 16:269–285. DOI: 10.1093/bioinformatics/16.3.269.

- Kotera M, Nishimura Y, Nakagawa Z, Muto A, Moriya Y, Okamoto S, Kawashima S, Katayama T, Tokimatsu T, Kanehisa M, Goto S. 2014. PIERO ontology for analysis of biochemical transformations: effective implementation of reaction information in the IUBMB enzyme list. *Journal of Bioinformatics and Computational Biology* 12:1442001. DOI: 10.1142/S0219720014420013.
- Maistrenko OM, Mende DR, Luetge M, Hildebrand F, Schmidt TSB, Li SS, Rodrigues JFM, von Mering C, Pedro Coelho L, Huerta-Cepas J, Sunagawa S, Bork P. 2020. Disentangling the impact of environmental and phylogenetic constraints on prokaryotic within-species diversity. *The ISME Journal* 14:1247–1259. DOI: 10.1038/s41396-020-0600-z.
- Marchesi JR, Ravel J. 2015. The vocabulary of microbiome research: a proposal. *Microbiome* 3:31. DOI: 10.1186/s40168-015-0094-5.
- McLaren MR, Willis AD, Callahan BJ. 2019. Consistent and correctable bias in metagenomic sequencing experiments. *eLife* 8. DOI: 10.7554/eLife.46923.
- Menzel P, Ng KL, Krogh A. 2016. Fast and sensitive taxonomic classification for metagenomics with Kaiju. *Nature Communications* 7:11257. DOI: 10.1038/ncomms11257.
- Miller CS, Baker BJ, Thomas BC, Singer SW, Banfield JF. 2011. EMIRGE: reconstruction of full-length ribosomal genes from microbial community short read sequencing data. *Genome Biology* 12:1–14. DOI: 10.1186/gb-2011-12-5-r44.
- Nasko DJ, Koren S, Phillippy AM, Treangen TJ. 2018. RefSeq database growth influences the accuracy of k-mer-based lowest common ancestor species identification. *Genome Biology* 19. DOI: 10.1186/s13059-018-1554-6.
- Pericard P, Dufresne Y, Couderc L, Blanquart S, Touzet H. 2018. MATAM: reconstruction of phylogenetic marker genes from short sequencing reads in metagenomes. *Bioinformatics* 34:585–591. DOI: 10.1093/bioinformatics/btx644.
- Puigbò P, Lobkovsky AE, Kristensen DM, Wolf YI, Koonin EV. 2014. Genomes in turmoil: quantification of genome dynamics in prokaryote supergenomes. *BMC Biology* 12. DOI: 10.1186/s12915-014-0066-4.
- Raymann K, Moeller AH, Goodman AL, Ochman H. 2017. Unexplored Archaeal Diversity in the Great Ape Gut Microbiome. *mSphere* 2. DOI: 10.1128/mSphere.00026-17.
- Sinha R, Abu-Ali G, Vogtmann E, Fodor AA, Ren B, Amir A, Schwager E, Crabtree J, Ma S, Abnet CC, Knight R, White O, Huttenhower C. 2017. Assessment of variation in microbial community amplicon sequencing by the Microbiome Quality Control (MBQC) project consortium. *Nature Biotechnology* 35:1077–1086. DOI: 10.1038/nbt.3981.
- Spencer SJ, Tamminen MV, Preheim SP, Guo MT, Briggs AW, Brito IL, A Weitz D, Pitkänen LK, Vigneault F, Juhani Virta MP, Alm EJ. 2016. Massively parallel sequencing of single cells by epicPCR links functional genes with phylogenetic markers. *The ISME journal* 10:427–436. DOI: 10.1038/ismej.2015.124.
- Thomas PD, Mi H, Lewis S. 2007. Ontology annotation: mapping genomic regions to biological function. *Current Opinion in Chemical Biology* 11:4–11. DOI: 10.1016/j.cbpa.2006.11.039.
- Willis AD. 2019. Rigorous Statistical Methods for Rigorous Microbiome Science. *mSystems* 4. DOI: 10.1128/mSystems.00117-19.

Reviewer 2

In this article, the authors propose an overview of the use of different approaches for microbiome data analyses, the questions that can be tackled using them, and their respective limitations. A particular focus is provided on the bioinformatics aspects, including an overview of the diversity of the most popular tools, and under which conditions/for which specific purposes they could be better used. Taxonomic and functional assignments tools are thoroughly discussed. And the crucial question of how to integrate the taxonomic and functional aspects. How marker-gene and shotgun metagenomic sequences (MGS) data are currently linked is exposed and the limitations of different approaches given. How the two approaches can lead to contradicting results, as well as the recurrent problem of reproducibility on microbiome data when using different bioinformatics pipelines are thoroughly discussed. Some interesting leads on how to make the field of microbiome biology more robust are given.

The paper is very well-written and thorough on several aspects, including by explaining the main trends in “DNA-based microbiome data” analyses. It is very interesting both from the level of technical details that are given, and from the fact that it does the synthesis of the current major pitfalls in microbiome studies. I particularly enjoyed the “Overview” and the last section on “Current state of the integration of taxonomic and functional data types”. As such, beyond proposing a view of the current state-of-the-art, I think this primer paper should contribute to the reflexion on what good practices could be taken, and which approaches are the most promising in order to make discoveries in microbiome studies more robust and reliable in the near future.

In general, I thought the titles of the big sections could be improved to better reflect their content. A few sections might require a bit of rewriting for clarification, and I would like to raise some points that are listed below.

1) The section “marker-gene sequencing”, where the case of the 16S rRNA amplicon sequencing is discussed at length (which is interesting!), is mostly dedicated to the particular task of characterizing the diversity within a community. However, it is only on lines 254-256 that the goal for using what is described as “robust marker genes” is introduced: “to characterize and compare the relative abundances of prokaryotes across communities.”

I think the 1st page of the manuscript could be re-arranged, and clarified to explain the particular usage of marker-gene approaches that is exemplified here.

- At the beginning of the section there is a discussion on the definition of a “robust marker gene”. But I believe this line of discussion depends on the goal of marker-gene sequencing – that should thus be introduced beforehand. Marker-gene approach can also be taken to question the presence of given metabolic processes in a particular environment. In which case, it is more important to fish for genes that are specifically involved in that process, leading even sometimes to multiply the set of probes to use in order to capture the diversity of the gene involved in the process of interest (some are paraphyletic for instance). In that case, the fact that the gene in question is a good molecular chronometer does not matter much, right? Or did I miss the point here?

- Line 156: a more general term would be “homolog”, as “ortholog” limits to vertically transmitted marker genes (excluding duplicated or laterally transferred genes for instance). Unless if it is explained beforehand that a desirable property of a marker gene could be to be vertically inherited? Or is the term “ortholog” used here to suggest a conserved biological function? Please clarify.

- In the end, I have the feeling that the first part of this 1st section kind of falls flat, as the authors write on lines 200-201: “Therefore, to select a robust marker gene one should adhere in some ways to the Goldilocks principle: some nucleotide conservation is needed, but not too much.” Maybe could this first part be shortened and be more straight-forward?

2) Lines 270-272: Please clarify what you mean by “V4-V5 region overrepresented Firmicutes ... while drastically underestimating Actinobacteria”. Do you mean that these regions are not present from Actinobacteria? Or that the diversity is over-estimated in Firmicutes and under-estimated in Actinobacteria based on this region? Same comment for line 290-291 for V1-V2 region.

3) In the section “Shotgun Metagenomics Sequencing”, I felt like the topic of the contribution of MGS approach and MAG (metagenome assembled genomes) reconstruction to explore extant biodiversity was somehow missing (CPR, DPANN, Asgard archaea...). MGS helped to reveal novelties both at the taxonomic and functional level. As a conceptual advantage of the MGS approach, in spite of some biases highlighted by the authors, is that it is not needed to have an a priori of what is looked for. This is how some entire clades of archaea were missed by 16S approaches because of the probes being designed from known diversity (e.g. Raymann et al 2017, mSphere).

- On lines 443-447 an example is given for taxa represented in 16S data but not MGS. To be fair, the converse is also true. I don't say the authors do not explicitly mention that there are caveats with both approaches, but this is one could be worth to be reminded.

- On lines 1004-1013, it could be added that techniques to bin MGS data as MAG could be a part of the solution.

4) On “the concordance of differential abundance results between actual and predicted metagenomics profiles” (lines 1882-1294), any lead on why the results are agreeing only “moderately well”?

5) Just a suggestion... Some figures could have been added to illustrate some parts of the text.

- On lines 1222-1225, the principle on which relies PICRUST for inferring function is introduced. It could have been illustrated by a figure.

- On lines 1409-1412, “stacked barplots” are mentioned to be used to study functional shifts. Such a typical plot could have been borrowed from a published study for instance?

6) In the Discussion part, it would have been interesting to have the authors opinions on the role that could play new sequencing techniques in the future to help with some of the issues presented? For instance, on the advent of long-reads sequencing for MGS? Don't you think it could eventually be a way to integrate taxonomic and functional analyses, by linking for instance 16S genes to big contigs, obtaining better quality MAGs, etc...?

7) Minor points and typos:

- A list of abbreviations should be included to help the reader. Otherwise, some of the less used abbreviations could be abandoned?

- Line 158: should it be “twice” instead of “double”?

- Line 1441 (and thereafter): maybe capitalize the tool name “phylogenize” to make it stand as a name in the text?

- Line 1445: “a taxa” => should be corrected by “a taxon”.

Reviewer 3

This review addresses many of the technical issues in the microbiome field. The text is very clear and concise, and it is very interesting for both initiated and uninitiated readers.

In general, the main point of the MS is the challenge of integrating taxonomic data with functional data. I agree that this is an issue but I feel in general the review downplay too much the binning/MAG approach dealing with this issue. I also missed in the text any discussion regarding long reads and how the 3rd generation sequencing methods could help with some of the limitations.

I have a few small comments that could improve the final version of the MS.

Line 107: There is often more statistical power to detect overall differences based on alpha and beta diversity metrics than to detect associations with individual features, but diversity-level insights are also less actionable (Shade 2017).

- However, often the difference of abundance in individual taxa/rank is larger than the difference in diversity indexes, especially in host-microbiome studies.

Line 422: This interest has culminated in the generation of enormous MGS datasets such as the ongoing work on the Earth Microbiome Project (Thompson et al. 2017) and the Human Microbiome Project (Lloyd-Price et al. 2017).

- Here another good and more recent example would be TARA oceans.

Line 548: “genes are expressed in cells, not in a homogenized cytoplasmic soup” (McMahon 2015).

- Agreed, however many ecological functions are performed in a collaborative way by consortiums.

Line 670: relative abundances by the mean relative abundance

- Should read geometric mean.

Line 723: This discussion of microbiome data characteristics has focused on taxonomic features based on either 16S sequencing or read-based MGS data analysis. However, it is important to emphasize that count tables produced from MAGs do not resolve this issue. In fact, attempting to account for these challenging characteristics of microbiome count data and the links between taxa and function makes the analysis more difficult.

- At the end of this, I would suggest a few lines about the network of co-abundances, for example using the SparCC tool.