



Peer Community Journal

Section: Genomics

RESEARCH ARTICLE

Published
2021-11-23

Cite as

Gavin M. Douglas and Morgan
G. I. Langille (2021) *A primer
and discussion on DNA-based
microbiome data and related
bioinformatics analyses*, Peer
Community Journal, 1: e5.

Correspondence

morgan.langille@dal.ca

Peer-review

Peer reviewed and
recommended by
PCI Genomics,

[https://doi.org/10.24072/pci.
genomics.100008](https://doi.org/10.24072/pci.genomics.100008)



This article is licensed
under the Creative Commons
Attribution 4.0 License.

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

Gavin M. Douglas ¹ and Morgan G. I. Langille ²

Volume 1 (2021), article e5

<https://doi.org/10.24072/pcjournal.2>

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.

¹Department of Microbiology and Immunology, McGill University, Montréal, Québec, Canada, ²Department of Pharmacology, Dalhousie University, Halifax, Nova Scotia, Canada

Peer Community Journal is a member of the
Centre Mersenne for Open Scientific Publishing
<http://www.centre-mersenne.org/>



CENTRE
MERSENNE

Contents

Background.....	2
Glossary	4
Phylogenetic marker-gene sequencing.....	5
Metagenomics sequencing.....	8
Characteristics of microbiome count data	10
Protein databases and ontologies for microbial genome functional annotation.....	13
Marker-gene-based metagenome prediction methods	18
Current state of the integration of taxonomic and functional data types	21
Outlook.....	25
Acknowledgements.....	28
Conflict of interest disclosure	28
References	28

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**. As the suffix "biome" suggests, a microbiome refers to these constituent elements in a defined habitat (Berg et al., 2020). 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, DNA sequencing has several advantages over culturing. First, 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 (Bukin et al., 2019; Jones et al., 2015; Watson et al., 2019). But nonetheless, DNA sequencing provides a more accurate view of the relative abundances of the community members than would be possible from culturing alone. Second, DNA sequencing is often a less time and labour-intensive method for assessing overall community diversity, although high-throughput culturing methods are becoming more common (Watterson et al., 2020). This is important, because high-throughput characterization of microbial communities is key to understanding microbial diversity, as closely related organisms can drastically differ in metabolic potential (Tettelin et al., 2005; Welch et al., 2002).

For these reasons, DNA sequencing remains the predominant method for characterizing microbial communities, although it is well-complemented by culturing (Lau et al., 2016). This method does have disadvantages however, for instance, it is difficult to distinguish between live, dormant, and dead cells using DNA sequencing (Carini et al., 2016). For researchers specifically interested in profiling active cells, leveraging alternative techniques such as metatranscriptomics, metaproteomics, and/or culturing, would be more appropriate.

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 sample environments (e.g., locations, disease states, etc.) and the relative abundance of features. A similar goal is often to investigate whether different measures of diversity in a studied dataset are associated with 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). In addition, many diversity metrics rely on unrealistic assumptions and there has been a recent push to develop more robust methods (Martin et al., 2020; Willis, 2019).

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 partially 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. Approaches that expand on basic marker-gene sequencing, such as epicPCR (Spencer et al., 2016), can provide information on the presence of small numbers of genes linked to marker-genes, but generally only limited genomic content can be gleaned from these methods. In contrast, **metagenomics sequencing** (MGS) provides information on all DNA present in a sample. MGS data can be used for analyzing both taxonomic and functional profiles. However, it is difficult to integrate the two data types, largely due to the complexity of microbial communities, the lack of robust databases, and the fragmented nature of DNA sequencing. In other words, it is relatively straightforward to identify genes in MGS data but challenging to determine from which genomes they originated.

Herein we introduce the key forms of these data types and highlight important caveats that should be considered when they are analyzed. Although many of our examples are taken from the human microbiome literature, our key points and suggestions are relevant to research in any microbial environment. We first cover the fundamentals of microbiome data analysis, starting with marker-gene sequencing, and then move to recently developed tools that could be leveraged to conduct joint analyses of taxonomic and functional data types. We conclude by highlighting two important challenges that must be addressed in microbiome data analysis.

Glossary

alpha diversity: The diversity in a single community (i.e., a particular sample). In the microbiome literature there are many alpha diversity metrics, such as richness, the Gini Simpson index, and the Shannon diversity index.

amplicon sequence variant: A single DNA sequence from a marker-gene sequencing dataset. These variants are produced through denoising methods, such as DADA2, deblur, and UNOISE3, which remove sequences that contain likely errors, rather than clustering them into operational taxonomic units. Due to this approach, amplicon sequence variants can in theory correspond to exact biological molecules and enable single-nucleotide differentiation between amplified sequences.

beta diversity: The diversity between communities (i.e., inter-sample distances). Similar to alpha diversity, many beta diversity metrics are applied in the microbiome literature, such as weighted and unweighted UniFrac distances, Bray-Curtis distance, and Aitchison distance.

contributinal diversity: The diversity of taxa that encode and/or perform a specific function. This is most commonly reported in terms of the diversity of prokaryotes with the potential to encode (or “contribute”) a specific pathway. Contributinal diversity can be reported in terms of either alpha or beta diversity. For alpha diversity comparisons, the Gini-Simpson index has been most commonly used.

function: A generic term with different definitions depending on the biology subdiscipline. In the microbiome literature a function is a general term referring to genome elements or biological processes performed by organisms. In practice, functions most commonly correspond to genes and pathways encoded by taxa. Genes are usually grouped into coarser categories known as gene families for microbiome analyses, which can be defined based on either sequence or functional similarity.

marker-gene: In the microbiome literature this most commonly refers to a gene that can be profiled to identify taxonomic lineages. We primarily discuss 16S rRNA gene sequencing as an example of marker-gene sequencing.

metagenome-assembled genomes: Genomes assembled from metagenomics sequencing data, without the requirement of culturing the organisms. The quality of these genomes is a contentious issue, as there is much higher risk of mis-assembling a genome in a mixed community compared with traditional genomes based on pure cultures.

metagenome prediction: Functional prediction of the genes and/or pathways present in a community based on taxonomic or marker-gene information. We primarily discuss PICRUSt2, a tool that we developed, which predicts the genome content for each query 16S rRNA gene through a phylogenetic approach.

metagenomics sequencing: (MGS) Untargeted sequencing of all DNA in a sample, which typically represents a mixture of organisms. This term is also sometimes used ambiguously to refer inclusively to both marker-gene and shotgun metagenomics sequencing, but herein it refers to only the latter. This data type requires more sophisticated processing and analyses than required for marker-gene investigations, but directly provides information on both the taxa and the genes encoded in the community. Ideally these data can be assembled into robust metagenome-assembled genomes. In practice the depth of sequencing is often insufficient to produce high-quality genomes and often instead

analyses focus on individual reads.

microbiome: A general term referring to a microbial community (and also the genes encoded, and metabolites produced) in a defined environment. The term microbiota is used when referring specifically to the taxa in a microbiome. In practice authors often use subtly different definitions of the term microbiome, such as referring to the bacterial portion of the community only, which has driven calls for more precise definitions in the field (Berg et al., 2020).

operational taxonomic unit: A pragmatic term used to denote a group of taxa as similar in some sense, typically for a particular study. In the microbiome field this term refers to clustering marker-gene sequences into operational groups. Several databases contain operational taxonomic unit sequences that can be re-used across studies. For 16S rRNA gene data analyses the traditional cut-off for defining operational taxonomic units is a sequence identity of 97%. This is a qualitatively different approach from denoising data to identify amplicon sequence variants.

pathway reconstruction: The processing of inferring whether a pathway is present or not based on the genes encoded in a set. For example, pathway reconstruction can be performed to infer which pathways could be potentially active given the genes encoded in a particular genome. In the microbiome field, this idea is commonly expanded so that gene family abundances from many taxa are used to infer which pathways could be active across the entire community.

Phylogenetic marker-gene sequencing

The earliest developed and most common form of microbiome sequencing is marker-gene sequencing, also known as amplicon sequencing. Under this approach specific genes are PCR-amplified and then sequenced. These genes can be markers of a particular functionality (e.g., Hug and Edwards, 2013), but more commonly this approach is employed to taxonomically profile a community, which is the purpose that we will discuss. Sequencing the 16S rRNA gene (hereafter referred to as 16S sequencing) is the most common amplicon sequencing approach for taxonomic profiling (see Box 1). Such 16S datasets are commonly produced to characterize and compare the relative abundances of prokaryotes across communities. However, despite the ubiquity of such datasets, they are non-trivial to analyze and interpret. There are numerous methodological reasons for this difficulty.

First, due to sequencing length constraints, only certain 16S rRNA gene variable regions are typically amplified and sequenced. Each variable region has particular strengths and limitations (Abellan-Schneyder et al., 2021; Chen et al., 2019; Johnson et al., 2019). Along with our colleagues we have previously compared the biases between the amplified fragments from variable regions four and five and from regions six to eight (written as V4-V5 and V6-V8, respectively) on a mock community from the Human Microbiome Project (Comeau et al., 2017). We found that sequencing the V4-V5 region resulted in a higher proportion of Firmicutes and Bacteroides and a lower proportion of Actinobacteria, compared with the known abundances. In contrast, sequencing the V6-V8 region resulted in a higher proportion of Proteobacteria and a lower proportion of Bacteroides. These biases highlight that choice of variable region can depend on which taxa are of interest. This is particularly true for less widely surveyed taxa such as archaea, which traditionally have been difficult to detect with 16S rRNA gene primers designed for bacteria (Bahram et al., 2019). For example, the V4-V5 region was recently shown to be superior to region V6-V8 for studying archaea in the North Atlantic Ocean (Willis et al., 2019). In this case the authors were particularly interested in archaeal diversity, so the V4-V5 region was more appropriate as it could be used to amplify the 16S rRNA gene of more archaea.

Box 1: Characteristics of robust marker-genes as exemplified by the 16S rRNA gene

There are two key requirements for robust marker-genes. First, they must be encoded by all (or at least most) taxa of interest. Second, the observed sequence divergence between homologs should be approximately equal to the neutral mutation fixation rate multiplied by double the divergence time between homologs (Woese, 1987). Note that the divergence time should be doubled because mutations could accumulate in either lineage since the organisms diverged. Genes displaying this second requirement have been referred to as molecular chronometers. This term highlights the close link between these marker-genes and the concept of the molecular clock (Zuckermandl and Pauling, 1965): given equal mutation rates and equal fixation rates for neutral mutations, the number of neutral substitutions between organisms is directly proportional to the evolutionary divergence between them.

However, there are many reasons why a gene might be an unreliable molecular chronometer (Janda and Abbott, 2007). One reason is that if a gene varies in function across taxa then contrasting selection pressures could result in different non-synonymous substitution rates (Wheeler et al., 2016). For instance, as previously observed (Woese, 1987), the cytochrome complex gene is a useful molecular chronometer in eukaryotes, but suffers from drawbacks. This gene was shown to be useful for building early phylogenetic trees that represented both long evolutionary distances across eukaryotes and short distances between human populations (Fitch and Margoliash, 1967). However, within prokaryotes the cytochrome complex systematically varies in size, which is believed to be due to positive selection (Ambler et al., 1979). Because positive selection is likely driving divergence between orthologous cytochrome complexes, in at least some cases it would be an invalid molecular chronometer to study in prokaryotes. Similarly, if a gene is sufficiently divergent between organisms then it can be difficult to accurately align residues. Misalignments lead to inaccurate estimates of evolutionary divergence, which is particularly true if the gene accumulates insertions and deletions. Such highly divergent regions, particularly in areas under no selective constraint, have been referred to as "evolutionary stopwatches" (Woese, 1987), because they are useful only at short evolutionary distances. 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.

The 16 Svedberg (16S) ribosomal RNA (rRNA) gene fits well with this principle. This gene features highly conserved regions surrounding nine less conserved regions (referred to as variable regions). It is also encoded by all prokaryotes and represents 50 helical RNA regions encoded by approximately 1,500 base-pairs (Woese et al., 1980). This high number of independent functional domains is valuable in a marker-gene (Woese, 1987). 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 evolutionary divergence. This gene also encodes a highly conserved function across both prokaryotes and eukaryotes (where it is called the 18S rRNA gene). The 16S rRNA molecule is part of the 30S small subunit of the ribosome, which helps initiate protein synthesis by binding the Shine-Dalgarno sequence in messenger RNA to align the ribosome with the encoded start codon. Many changes in the highly conserved regions of the 16S rRNA gene affect its binding affinity to the ribosome and messenger RNA. The strong negative selection acting against such substitutions makes these regions valuable for detecting rare substitutions, anchoring alignments, and for primer design (Wang and Qian, 2014).

Since the 16S rRNA gene was identified as a useful molecular chronometer, it has been the prime marker-gene used to develop phylogenetic models of the tree of life. Most famously, an alignment of 16S (and 18S) rRNA gene sequences from across life lead to distinguishing archaea, bacteria, and eukaryotes into distinct domains (Woese and Fox, 1977). In these early days, research focused on analyzing the rRNA sequences of isolated microbes. This was painstaking work, as illustrated by the prediction in 1987 that future research groups could plausibly sequence on the order of one hundred 16S rRNAs a year (Woese, 1987).

Thirty-four years later, through next-generation sequencing technology, insufficient availability of sequenced rRNA genes is no longer a common complaint. Databases such as SILVA contain enormous collections of sequenced small subunit fragments; as of August 2020 SILVA contained 9,469,124 non-clustered, independent sequences (Quast et al., 2013).

Typically, however, the taxonomic scope of interest and region biases in a particular environment are not clear and little or no rationale is given for the variable region selection. This is a problem, because analyses of the same communities with different variable regions can result in not only systematic biases in the raw data, but also in strikingly different biological interpretations. For example, key species that modulate human vaginal health are underrepresented or missing in V1-V2 sequencing datasets, such as *Gardnerella vaginalis*, *Bifidobacterium bifidum*, and *Chlamydia trachomatis* (Graspeuntner et al., 2018). Application of this region for profiling vaginal samples, instead of the more appropriate choice of the V3-V4 region, can result in entirely missing associations between vaginal health and the microbiome. Similarly, a comparison of the tick microbiome based on six sequenced 16S rRNA gene regions found a wide range of the number of prokaryotic families and in the Shannon diversity index for each individual tick (Sperling et al., 2017). The problem of such biases in variable region selection is beginning to recede as long-read technologies, such as that developed by Pacific Biosciences of California and Oxford Nanopore Technologies Limited, enable full-length 16S sequencing (Callahan et al., 2019; Johnson et al., 2019). However, it will remain an important issue for the foreseeable future as long as the microbiome is largely studied by short-read sequencing.

Regardless of the sequenced region, most reads originating from the same biological molecule will differ due to sequencing errors. Raw reads are either clustered based on sequence identity into **operational taxonomic units** or alternatively errors are corrected to produce **amplicon sequence variants**. Operational taxonomic units are typically clustered at 97-99% identity (Goodrich et al., 2014), which often results in merging different species into a single operational taxonomic unit (Mysara et al., 2017). This issue has long plagued 16S rRNA gene-based analyses. For instance, *Bacillus globisporus* and *Bacillus psychrophilus* are problematic cases because their 16S genes share 99.5% sequence identity, but are highly distinct at the genome level (Fox et al., 1992).

In contrast to clustering approaches, error-correcting approaches, referred to as denoising methods, theoretically can correct raw reads sufficiently well to produce exact biological molecules. Several different denoising approaches have recently emerged. DADA2 is the most sophisticated approach, which generates a different parametric error model for every input sequencing dataset (Callahan, McMurdie, et al., 2016). The raw sequencing reads are then corrected to generate amplicon sequence variants based on this error model. Deblur (Amir et al., 2017) and UNOISE3 (Edgar, 2016) are two other denoising tools that are based on rapidly clustering raw reads and using predetermined hard cut-offs related to the expected error rates to generate amplicon sequence variants. We and other colleagues have evaluated the performance of these three tools and open-reference operational taxonomic unit clustering (which combines both de novo and reference-based clustering) and found that all three denoising methods result in similar overall microbial communities (Nearing et al., 2018). In contrast, we found that open-reference operational taxonomic unit clustering resulted in a high rate of spurious identifications compared to these methods. Nonetheless, there were important differences between the three denoising methods, particularly in terms of richness and when profiling rare taxa (Nearing et al., 2018). A more recent independent validation based on a higher number of test datasets reached similar conclusions (Prodan et al., 2020).

In addition to 16S rRNA gene sequencing data, there are multiple marker-genes appropriate for profiling eukaryotic diversity. The 18S rRNA gene is the homolog of the 16S rRNA gene in eukaryotes and is widely used to profile that domain. However, fungi are more difficult to distinguish based on the 18S rRNA gene, because fungi lack several variable regions for this gene (Schoch et al., 2012). Instead, the internal transcribed spacer region, although not strictly a marker-gene, is more often amplified to study fungal communities, because it typically has more resolution to distinguish fungi than the 18S rRNA gene (Liu et al., 2015). This region is within the nuclear rRNA cistron of fungi genomes, which contains the 18S, 5.8S, and the 28S rRNA genes. The internal transcribed spacer regions encompasses the two intergenic regions, which have relatively high rates of insertions and deletions, and the 5.8S rRNA gene (Schoch et al., 2012). Only a single intergenic region is typically amplified, referred to as regions one and two,

which have better discriminatory resolution for the major phyla Basidiomycota and Ascomycota, respectively (Bellemain et al., 2010).

Although the marker-genes described above are the most commonly profiled loci, in many cases there are marker-genes more appropriate for specific lineages. For example, several halophilic species of *Haloarcula* encode multiple 16S copies that can differ by more than 5% sequence identity within the same genome (Sun et al., 2013). Consequently, different marker-genes are often used when building phylogenetic trees representing a single species or genus.

The chaperonin-60 (*cpn60*) gene is one useful alternative prokaryotic marker-gene, which is particularly useful for distinguishing taxa at resolutions below the genus level (Links et al., 2012). For example, the *cpn60* gene has been frequently profiled in vaginal microbiome samples, because variation at this locus can distinguish subgroups of *Gardnerella vaginalis* that cannot be distinguished based on the 16S rRNA gene alone (Jayaprakash et al., 2012). Similarly, the gene *rpoB*, which encodes the DNA-directed RNA polymerase subunit beta, is another valuable prokaryotic marker-gene, which provides comparable or better taxonomic resolution to the 16S rRNA gene (Case et al., 2007). Profiling *rpoB* can sometimes better identify relevant taxa in a community. For instance, it has been used for identifying a known nematode symbiont missed by standard 16S profiling with the V3-V4 region (Ogier et al., 2019).

More generally, marker-genes for specialized comparisons are often chosen to match the defining function of a given lineage. For example, the methyl coenzyme M reductase A 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.

Metagenomics sequencing

Metagenomics sequencing (MGS) is a qualitatively different method from marker-gene sequencing, because it involves sequencing all DNA in a community. This is a major advantage and means that MGS can profile any DNA-encoding taxa, including DNA viruses and microbial eukaryotes. This has enabled the discovery of novel lineages, including previously unknown phyla (Spang et al., 2015), through analyzing MGS data. However, this characteristic also makes data analysis more challenging. This is particularly because sources of DNA that are not of interest, such as host DNA or contaminants (especially in low biomass samples), can often be substantial proportions of MGS datasets. 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 investigation 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 Earth Microbiome Project (Thompson et al., 2017), the Human Microbiome Project (Lloyd-Price et al., 2017b), and the TARA Oceans investigations (Sunagawa et al., 2015).

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 metagenomics 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 particular result is likely dependent on the taxonomic profiling approach used (see below).

Nonetheless, this wide variation in results highlights that any interpretation of MGS profiles, similar to 16S profiles, should be done cautiously. It is crucial to appreciate that any approach will have important weaknesses and that the generated profile will only partially represent the actual microbial diversity.

With those important caveats in mind, an understanding of the different approaches is nonetheless important to give context to MGS data analysis. Read-based workflows involve little or no assembly of the reads and instead each read (or pair of reads) is treated independently. This is the most common approach for analyzing MGS data, particularly because it can be performed with low sequencing depth (Hillmann et al., 2018) and in complex communities (Zhou et al., 2015). However, an important disadvantage of this approach is that taxonomic and functional annotations are typically generated and treated as entirely independent data types (Figure 1a). It is also possible to map reads against a set of known reference genomes, which does link the two data types (Figure 1b). Although this is an invaluable approach when applied to genomes assembled from the study environment (see below), the results are typically near incomprehensible when reads are mapped against a database of thousands of genomes at the nucleotide level. In contrast, when reads are mapped against reference genomes in protein space, using a tool such as Kaiju (Menzel et al., 2016), this approach can provide useful taxonomic profiles. Nonetheless, the most common approaches for generating taxonomic profiles are based on either a marker-gene or k-mer method.

Marker-gene approaches are based on the insight that specific genes can be used to identify the presence and relative abundance of certain taxa. An extreme example is to use solely the 16S rRNA gene for taxonomic classification (Hao and Chen, 2012). Several methods have been developed specifically for targeted assembly of this and other rRNA genes from MGS data (Gruber-Vodicka et al., 2020; Miller et al., 2011; Pericard et al., 2018). More commonly, marker-gene approaches base classifications on many genes. For instance, PhyloSift (Darling et al., 2014) leverages 37 nearly universal prokaryotic marker-genes (Wu et al., 2013) in addition to eukaryotic and viral gene sets to make a combined set of approximately 800 (mainly viral) gene families for classification. Aligned reads are placed into a phylogenetic tree of reference sequences and taxonomic classification is performed based on summing the likelihood of each taxa based on each read placement (Darling et al., 2014). MetaPhlAn is a contrasting approach that instead bases taxonomic predictions on the presence of clade-specific marker-genes, which are genes only found in that given lineage, and found in all members (Truong et al., 2015). This method has rapidly become the most popular marker-gene MGS approach. However, given that this approach is limited by the existence of robust clade-specific genes, it is not surprising that it tends to have low sensitivity (Miossec et al., 2020; Tessler et al., 2017), meaning that it misses taxa that are actually present.

In contrast, k-mer-based approaches are much more sensitive but have slightly lower specificity than marker-gene methods (Miossec et al., 2020). These approaches search for exact matches of short DNA sequences (k-mers) within reference genomes. An algorithm such as lowest-common ancestor is then performed to determine the likely taxonomic classification based on all matching genomes. Two common kmer-based approaches are kraken2 (Wood et al., 2019) and centrifuge (Kim et al., 2016), both of which match k-mers against a compressed database of reference genomes. One disadvantage of such methods is that taxonomic classifications can be highly dependent on the size of the database used (Nasko et al., 2018). In addition, the main challenge of analyzing taxonomic profiles output by these methods is the high number of rare taxa of different ranks identified, some of which may be false positives. Summarizing the output profiles with an additional approach, such as the Bayesian abundance re-estimation tool Bracken (Lu et al., 2017) in the case kraken2 data, can help mitigate these problems.

Most functional read-based methods are based on a similarity search of reads against a database of known gene families. This is primarily done in protein space, because protein similarity matches are more informative and the database requirements are lower (Koonin and Galperin, 2003). The common similarity searching tool BLASTX is prohibitively slow when scanning millions of reads, which has driven the development of faster alternatives like DIAMOND (Buchfink et al., 2015) and MMseqs2 (Steinegger and Söding, 2017). These faster alternatives are leveraged

by workflows implemented in software such as MEGAN (Huson et al., 2007) and HUMAnN2 (Franzosa et al., 2018) to identify gene family matches and output overall metagenome profiles. HUMAnN2 is a unique approach in that it first screens reads that map to reference genomes of taxa identified as present with MetaPhlAn2. This step enables a small subset of gene families to be linked directly to particular taxa. However, the vast majority of gene families typically have no taxonomic links and are only part of the community-wide metagenome. There are clear issues with the general approach implemented by these gene profiling approaches, as has been previously observed: "genes are expressed in cells, not in a homogenized cytoplasmic soup" (McMahon, 2015).

Linking functional annotations to specific taxa by assembling raw reads is the ideal approach to resolve this problem, but this too comes with caveats. Most importantly, insufficiently high read depth, which depends on the complexity of a sample, can result in too few assembled contigs to sensibly analyze. Nonetheless, with sufficiently high read depth metagenome assembly can be a valuable way to leverage information about microbial communities (Figure 1c). There are many metagenome assembly tools available, such as MetaSPAdes (Nurk et al., 2017) and MEGAHIT (Li et al., 2015). The resulting assembled contigs from these approaches are typically categorized (or "binned") into groups of contigs with similar characteristics. This binning is primarily performed by identifying contigs that are found at similar relative abundances across samples and/or that contain similar proportions of different k-mers (Ayling et al., 2020). These bins represent **metagenome-assembled genomes** that must undergo stringent checks to help evaluate the overall quality (Bowers et al., 2017). The key method for performing quality control on these genomes is to scan for known universal single-copy genes, with a tool such as CheckM (Parks et al., 2015). The percentage of universal single-copy genes present provides an estimate of overall genome completeness. In contrast, the number of universal single-copy genes found in multiple copies can be used to calculate the redundancy, which is potential evidence for contamination or strain heterogeneity in the genome. For further details on metagenomics assembly and binning tools, readers can find recent reviews that describe the available bioinformatics tools (Ayling et al., 2020; Breitwieser et al., 2019).

One final consideration is that several recent technologies have been developed that can lead to higher quality metagenome-assembled genomes. These include long-read sequencing technology (McCarthy, 2010; Mikheyev and Tin, 2014), Hi-C sequencing (Belton et al., 2012), optical mapping (Hastie et al., 2013), read clouds (Bishara et al., 2018), and single-cell metagenomics (Xu and Zhao, 2018). We have previously discussed the utility of these specific technologies in the context of producing improved metagenome-assembled genomes in more detail (Douglas and Langille, 2019).

Characteristics of microbiome count data

Regardless of the sequencing technology and workflow used for taxonomically profiling a microbial community, the final product is typically an abundance table. This is true for many sequencing approaches, such as in transcriptome datasets, but there are several important differences. First, unlike in the case of transcriptome read count tables where there are typically a known number of genomic loci, novel taxa and functions are frequently identified in microbiome data. For instance, novel operational taxonomic units, amplicon sequence variants, and contigs are frequently identified in taxonomic analyses. Similarly, 25-85% of proteins in MGS are novel microbial genes of unknown function (Prakash and Taylor, 2012). Second, no statistical distribution fits microbiome data in all contexts. For example, many statistical distributions, including the negative binomial (Love et al., 2014), beta binomial (Martin et al., 2020), and Poisson (Faust et al., 2012) distributions have been proposed as appropriate fits to microbiome data. However, upon analysis with real data these and other distributions fit with inconsistent accuracy (Calgario et al., 2020; Weiss et al., 2017). Last, microbiome abundance tables typically have high sparsity, meaning that there is a high proportion of features not found across many samples (Thorsen et al., 2016). These characteristics make microbiome data analysis challenging for all taxonomic analyses and most functional analyses.

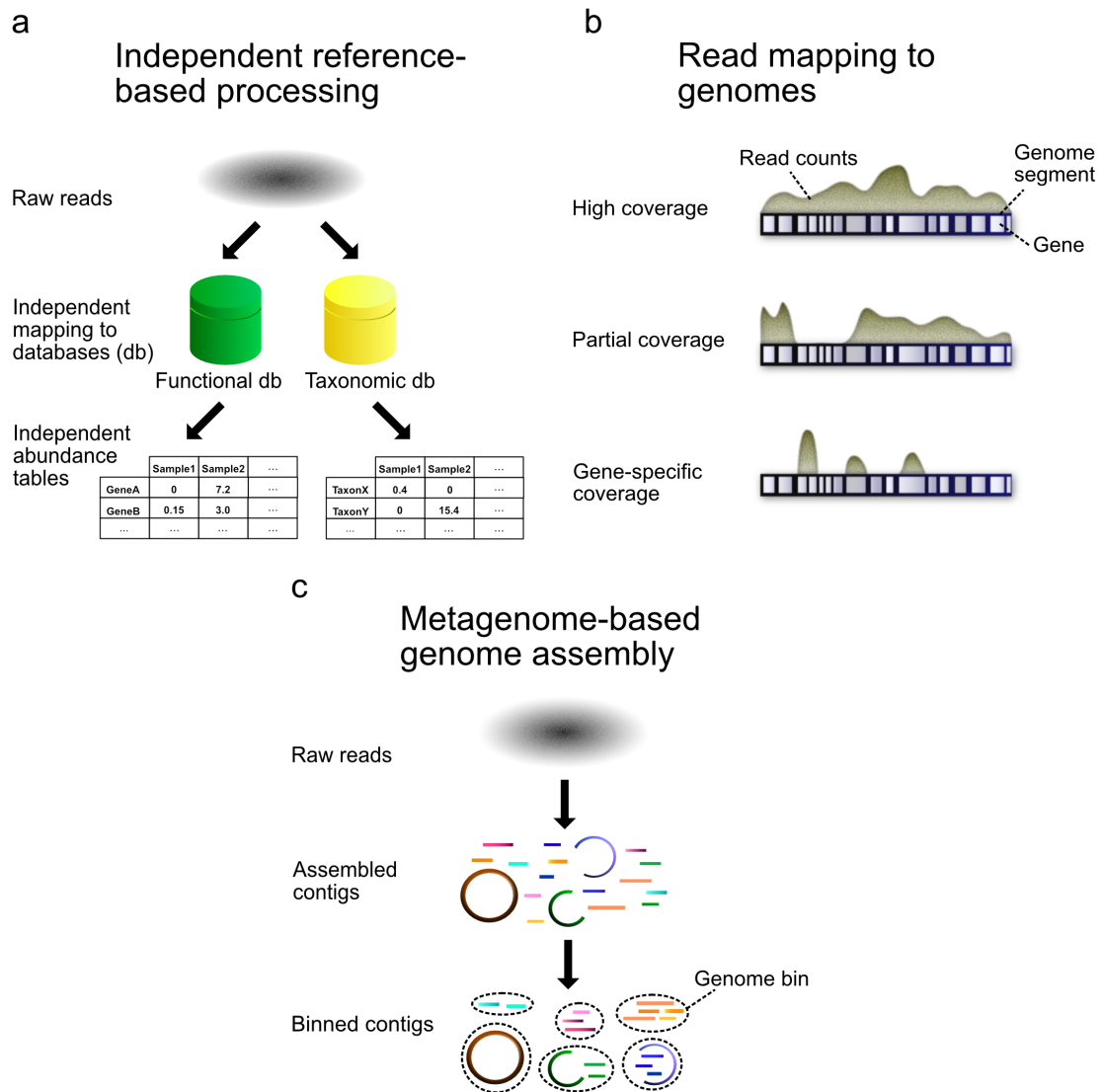


Figure 1: Key approaches for generating joint taxonomic and functional data from metagenomics sequencing data. (a) Read-based processing of metagenomics data to generate functional and taxonomic abundance tables independently. (b) Read mapping to genome sequences can be used to infer the presence of 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 community complexity is low.

These challenges are exacerbated by the inherent compositionality of sequencing data. Compositional data refers to data that is constrained to an arbitrary constant sum (Aitchison, 1982), such as the arbitrary number of raw sequencing reads output per sample. This characteristic

means that the observed abundance of any given feature is dependent on the observed abundance of all other features (see Figure 2 for an illustrative example of the implications of this characteristic). This implies a necessary consideration regarding microbiome sequencing data analysis: it only provides information on the relative abundances, or percentages, of features and does not provide insight on feature absolute abundances.

This important characteristic was not widely appreciated in the field until relatively recently, when researchers identified fatal issues with common approaches for analyzing microbiome data (Gloor et al., 2016, 2017). Standard differential abundance approaches, such as the t-test and Wilcoxon test, when applied to relative abundances, and microbiome-specific tools such as LEfSe (Segata et al., 2011) do not account for this compositionality. Common summary metrics for microbiome data, such as the UniFrac distance, also suffer from this problem (Gloor et al., 2017). This is a major issue, because ignoring this characteristic is known to lead to spurious discoveries with compositional data (Aitchison, 1982; Fernandes et al., 2014; Jackson, 1997).

Fortunately, there is active work in the field to resolve this issue and numerous compositional approaches have been developed. For instance, several compositional correlation approaches are now available (Friedman and Alm, 2012; Kurtz et al., 2015; Schwager et al., 2017). One such approach is SparCC, which computes inter-taxon correlations while accounting for artifactual correlations that occur simply due to the interdependency between features in the same compositional dataset (Friedman and Alm, 2012). Differential abundance approaches appropriate for compositional data analysis have also been developed, such as ALDEx2 (Fernandes et al., 2014, 2013) and ANCOM (Mandal et al., 2015). A common theme of these compositional approaches is that the data is transformed based on the ratio of feature relative abundances to some reference frame (Aitchison, 1982; Morton et al., 2019). This choice of reference frame varies substantially between approaches. For instance, ALDEx2 transforms relative abundances by the centred log-ratio transformation (Fernandes et al., 2013), which essentially normalizes feature relative abundances by the geometric mean relative abundance per sample. This approach transforms the original data but maintains the interpretation of individual features. In contrast, it has been suggested that analyses could instead be based on ratios between features (Morton et al., 2019), which converts the data type into comparisons of features rather than individual features.

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

To compound these discrepancies, there are even disagreements regarding how to preprocess and filter datasets prior to statistical testing. For instance, microbial features with low prevalence or that are only found at low read depths are often discarded. Ad-hoc cut-offs for feature filtering, such as a minimum prevalence of 10%, are often used, but there is little consistency across studies. In addition, it has been suggested that filtering out rare features based on read depth can, at least under certain conditions, reduce statistical power (Schloss, 2020).

Given the wide variation in differential abundance tool performance and unclear best-practices, how is a microbiome researcher to proceed? One possible answer is that a change in expectations regarding the interpretability of microbiome data analysis is needed. In particular, analyses using ratios between the relative abundances of taxa have been shown to be robust, although the increased robustness comes at the cost of interpretability (Morton et al., 2019). However, an important issue is how to determine which taxa should be the numerator and denominator

of each ratio. One solution is to leverage the bifurcating structure of a clustered tree (Morton et al., 2017; Pawlowsky-Glahn and Egozcue, 2011) or phylogenetic tree (Silverman et al., 2017) of features. Analyses can be focused on the ratios in relative abundances between features on the left-hand and right-hand of each node in the tree. Despite the potential of this approach, it is rarely used for standard microbiome analyses because it is unclear how to biologically interpret any differences in the values of these ratios across samples.

An alternative solution is to leverage additional data to transform relative abundances to absolute abundances. This alternative data could be quantitative PCR data, flow cytometry data (Vandeputte et al., 2017) or spiked-in sequences of known abundances (Zemb et al., 2020). Different sample preparation protocols prior to DNA sequencing can also help retain information about differential absolute amounts of DNA across samples as well (Cruz et al., 2021). These are exciting approaches, but they have not been validated across many datasets and at the moment there is no consensus regarding which methods perform best.

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 metagenome-assembled genomes 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.

Protein databases and ontologies for microbial genome functional annotation

To this point we have only discussed functional microbiome data in vague terms as referring to microbial gene abundances. When based on DNA sequencing data this information summarizes the functional potential, meaning the functions that are present, but not necessarily active in a community. However, rather than individual gene sequences, research is typically focused on gene families, which are defined based on close sequence identity and/or similar functionality from the gene's eye view. Alternatively, the focus is sometimes on higher-order functional categories like pathways, which represent functionality of groups of potentially interacting gene families in reactions. To complicate matters further, there are several different functional databases and ontologies for annotating microbial functions. Ontologies are representations of information groupings and relationships of arbitrary entities (Thomas et al., 2007). In the context of functional annotation, different ontologies represent different ways of functionally grouping genes and also of defining higher-level and more general microbial functions. Depending on which of these functional ontologies and sub-categories are analyzed, the characteristics of the data can drastically differ.

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

The UniRef database contrasts with another common functional ontology, the Kyoto Encyclopedia of Genes and Genomes (KEGG) database (Kanehisa et al., 2016; Ogata et al., 1999). KEGG is based on 23,530 individual gene families (as of September 10th, 2020), which are called KEGG orthologs. The advantage of KEGG orthologs is that the majority have well-described molecular functions that can be linked to higher-order KEGG pathways and modules. Accordingly, any analysis of KEGG data will likely result in less sparse count tables than the corresponding UniRef-based database, simply because KEGG orthologs are shared across more taxa than UniRef clusters.

To illustrate this point, we and our colleagues have previously compared the taxonomic coverage of each function within these two functional ontologies and each sub-category (Inkpen

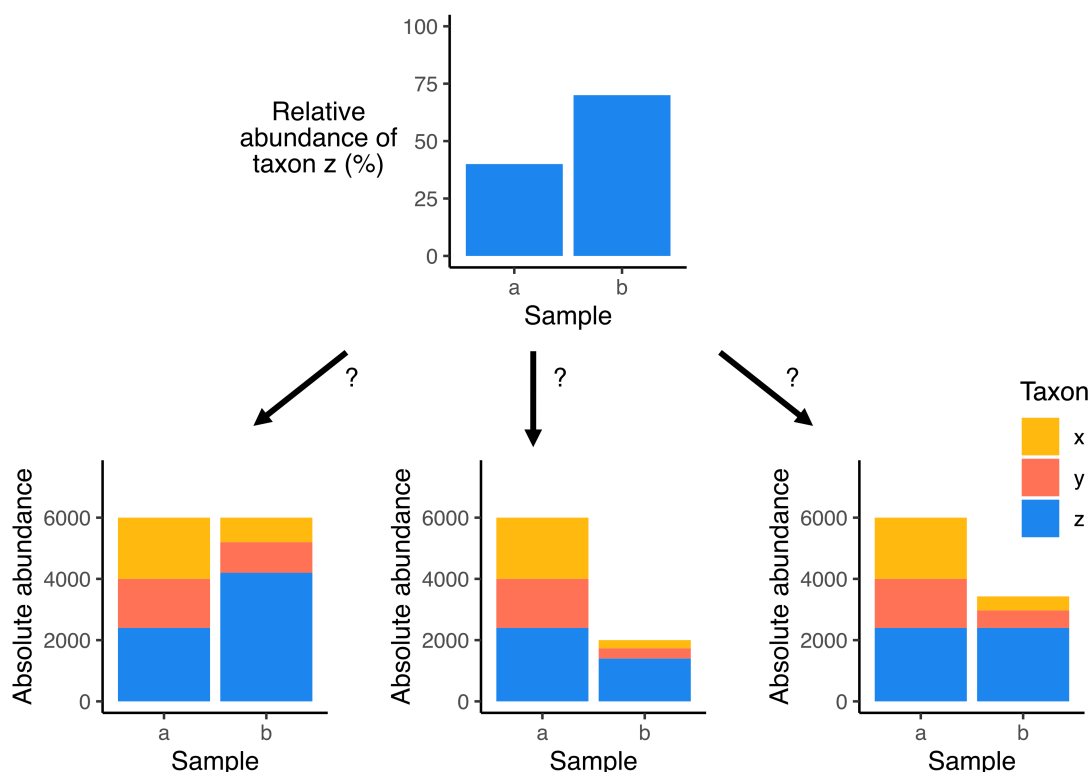


Figure 2: Microbiome sequencing data provides information only on relative abundances. This is because the data are compositional, meaning that it is constrained to sum to some arbitrary number, rather than corresponding to actual absolute abundances (e.g., cell counts or colony forming units). This illustrated example highlights how a difference in relative abundance in a given taxon, taxon z, between two samples does not provide information on whether there is a difference in terms of the absolute abundance of this taxon (except for certain circumstances noted in the main-text). From bottom left to right three possible configurations of absolute abundance are shown. The left panel shows the case where the total abundance of taxa is the same in each sample and so taxon z is indeed at higher absolute abundance in sample b. However, it is also possible that the absolute abundance of this taxon could be lower (bottom centre) or the same (bottom right) depending on the total absolute abundance in each sample.

et al., 2017). We found that all UniRef functions, including those in UniRef50 clusters, are on average found in a single domain and encoded by fewer than four species. In contrast, we found that KEGG orthologs were encoded in 1.3 domains and 184.3 species on average. Similarly, the high-level KEGG modules and pathways were predicted to be potentially active in a mean of 1.7 and 2.5 domains and 671 and 1267.6 species, respectively (Inkpen et al., 2017). Based on these statistics, clearly a shift in the abundance of a UniRef cluster should not be treated the same as a KEGG function: the former corresponds to the activity of a small number of species while the latter could correspond to a large assemblage. This example highlights that the choice of how function is defined in a given analysis can have profound effects on the biological interpretation.

In addition to UniRef and KEGG, several other functional ontologies have been leveraged for microbiome analyses. Key examples of additional function types include: Clusters of Orthologous Genes (COGs) (Makarova et al., 2015; Tatusov et al., 2000), Enzyme Commission numbers, Protein families (Pfam) (Finn et al., 2014; Punta et al., 2012), and TIGRFAMs (Haft et al., 2003). These categories represent a range of approaches for defining gene families and functional categories.

The COG strategy for functional annotation was originally intended to phylogenetically classify proteins into groups of orthologs (Tatusov et al., 2000). This one-to-one approach of matching individual orthologs has now been expanded to allow for more complex relationships between genes, such as paralogs and horizontally transferred homologs (Galperin et al., 2019; Makarova et al., 2015). As of 2015, there were 4,631 independent COGs (Galperin et al., 2015). The COG framework is similar to that of the eggNOG database (Jensen et al., 2008), which is a more high-throughput, automated approach. However, the key advantage of the COG database is that orthologous genes are clustered into 26 interpretable functional categories, which are expanded from categories originally defined to functionally bin *Escherichia coli* genes (Riley, 1993).

The Enzyme Commission number framework, which was developed in 1992 by the "International Union of Biochemistry and Molecular Biology", is a contrasting approach for functional annotation. Instead of focusing on orthologous genes, Enzyme Commission numbers specify particular enzyme-catalyzed reactions. An interesting characteristic of this database is that these reactions can be performed by non-homologous isofunctional enzymes (Omelchenko et al., 2010). As of August 12th, 2020, there were 6,520 Enzyme Commission numbers, which correspond to one of four levels of granularity. For example, the accession 3.5.1.2 corresponds to glutaminases, while the higher-level categories correspond to hydrolases (3.-.-), that act on carbon-nitrogen bonds other than peptide bonds (3.5.-), and that are in linear amides (3.5.1.-). One major advantage of Enzyme Commission numbers is that because they specify exact enzymatic reactions they are straight-forward to link into pathway ontologies based on reactions, such as MetaCyc pathways (Caspi et al., 2013).

The Pfam database categorizes protein families, which are protein regions that share sequence homology (Punta et al., 2012). Individual proteins with multiple domains can thus belong to multiple Pfam families. Each Pfam family is represented by a hidden Markov model, which models the likely amino acids at each residue and the likely adjacent amino acids based on curated alignments of representative protein sequences. This approach identified homologous protein regions, which are often hypothesized to have a shared evolutionary history, but not necessarily. As of May 2020, there were 18,259 Pfam families.

Lastly, TIGRFAMs are manually curated protein families, which are also identified based on hidden Markov models, but also additional pertinent information (Haft et al., 2003). As of September 16th, 2014, there were 4,488 TIGRFAMs. The distinguishing feature for this database is that different information supplements each hidden Markov model. For instance, certain TIGRFAM are annotated based on species metabolic context and neighbouring genes, while others are based on validated functions from the scientific literature. This database has been less commonly analyzed in recent years and is best known as the annotation system for early large-scale metagenomics projects (Venter et al., 2004). Alternative approaches, such as the FIGfam protein database are now more commonly used than TIGRFAMs. FIGfams are based on a similar approach, but instead of being manually curated they are aggregated into isofunctional groups based on shared roles in specific subsystems (Meyer et al., 2009).

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

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

However, it is worth noting that pathways identified based on DNA sequencing are merely theoretical reconstruction based on the identified individual gene families. Although there are several [pathway reconstruction](#) approaches, they all require some mapping from gene families or reactions to pathways. This mapping can be structured, meaning that optional and required contributors can be specified, or non-structured, meaning that all genes and/or reactions are treated equally.

The naïve approach for pathway reconstruction is to assume that a pathway is present if any gene or reaction involved is present in the community. This was the predominant approach used for pathway inference in early functional analyses (Meyer et al., 2008; Moriya et al., 2007) and in several pathway inference tools such as PICRUSt (Langille et al., 2013). Pathway abundance under this framework is calculated by summing the abundance of each contributing gene family. This approach errs towards avoiding missing the presence of a pathway, which is a concern in metagenomes as key genes may be missing due to mis-annotations. However, this approach comes at the cost of spurious annotations. Based on the naïve mapping approach the human genome was previously annotated as including the KEGG pathway equivalent of the reductive carboxylate cycle (Ye and Doak, 2011). This pathway is restricted to autotrophic microbes and is similar to reversing the Krebs cycle. Consequently, several gene families are shared in both processes. Under the naïve mapping approach, the presence of genes involved in the Krebs cycle are also evidence for the predicted presence of this atypical microbial pathway in humans. Similarly, vitamin C biosynthesis would also be predicted in humans based on the naïve approach (Ye and Doak, 2011). However, the *GLO* gene, which encodes the protein involved in the key last step of vitamin C biosynthesis in mammals, is pseudogenized in humans (Drouin et al., 2011), which makes vitamin C biosynthesis impossible.

The Minimal set of Pathways (MinPath) approach is an approach developed to address this issue (Ye and Doak, 2011). This tool identifies the smallest set of pathways, based on maximum parsimony, that are required to explain the presence of a set of proteins. In this way, the approach is more conservative than naïve mapping and also accounts for incomplete protein sets. This method has been applied in numerous contexts, including for the "Human Microbiome Project Unified Metabolic Analysis Network 2" (HUMAN2) (Abubucker et al., 2012; Franzosa et al., 2018) MGS gene family profiling and pathway reconstruction framework. This popular framework reconstructs pathways based on MinPath and infers pathway abundance based on different approaches, depending if the pathway mapping is structured. For unstructured mappings, the arithmetic mean of the upper half of individual gene family abundances is taken to be the pathway abundance (Abubucker et al., 2012). For structured mappings, the harmonic mean of the key (i.e., required) genes families is computed for pathway abundance (Franzosa et al., 2018). Both these approaches are motivated by the need to be robust to variable abundance in alternative gene families.

Although this approach for MGS pathway reconstruction is commonly performed, it is important to emphasize that it has not been universally accepted and there remains disagreement about best-practices. For example, "Evidence-based Metagenomic Pathway Assignment using geNe Abundance DATA" (EMPANADA) is a method that addresses the same issue as MinPath and HUMAN2 in a different way (Manor and Borenstein, 2017a). Pathway reconstructions from this tool are based on the distinction between genes that are shared with multiple pathways from those that are unique to a single pathway. Pathway support weightings are first given by the average abundance of gene families unique to each given pathway. The abundance of all shared gene families is then partitioned between all pathways according to their relative support values. Pathway abundances are then taken as the sum of the unique gene family relative abundances and the partitioned relative abundances of the shared gene families (Manor and Borenstein, 2017a).

The exact reconstructed pathways and their respective abundances differ depending on whether naïve mapping, MinPath/HUMAN2, or EMPANADA are used. Validating pathway reconstructions is challenging without a gold-standard comparison, particularly in metagenomes. Even in isolated genomes, as demonstrated by the above examples of the human pathway reconstructions, pathway reconstruction is non-trivial. However, the advantage in these cases is

Box 2: Comparing taxonomic and functional stability across communities

An introduction to microbiome functional data would be incomplete without addressing its ostensible high stability. Functional pathways are commonly at similar relative abundances across the same sample-types whereas taxonomic features, such as phyla, can substantially vary (Burke et al., 2011; HMP-consortium, 2013; Louca et al., 2016; Turnbaugh et al., 2009). This functional consistency is often taken to be evidence of environmental selection for particular microbial functions (Louca and Doebeli, 2017; Turnbaugh et al., 2009). However, the validity of comparing variation between these two data types is rarely discussed. We and our colleagues investigated this question from a philosophical perspective and concluded that any meaningful comparison of the relative variation between taxonomic and functional profiles is likely impossible (Inkpen et al., 2017). This difficulty is largely because it is unclear which levels of granularity would be meaningful to compare between each data type. For instance, the gene and pathway perspectives of function represent two extremes of functional granularity. Many different functional ontologies exist as well for defining functional groups, as discussed in the main-text. Because taxa and functional data types are qualitatively different from each other, the choice of how to compare the two is based on somewhat arbitrary decisions on how to categorize them.

This can be illustrated by comparing taxa and functions in the same communities based on different groupings of each data type. As described in the main-text, 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, Polz, et al., 2018). Nonetheless, this example demonstrates once again a common theme throughout this work: "function" has many meanings.

Box 3: DNA hybridization and early 16S rRNA gene studies established high genomic variability

Classic DNA hybridization experiments highlighted the high genomic variability between different bacteria (Brenner, 1973; Mandel, 1966). These experiments were based on mixing single-stranded DNA from two organisms and recording the melting temperature required to separate the strands. Higher melting temperatures are required to break apart DNA that shares more complementary bases connected by hydrogen bonds. Accordingly, this approach provides a rough estimate of the genetic distance between different strains or species.

An early comparison of these genetic distances with 16S dissimilarity across 34 bacteria computed a linear correlation of 0.728 (Devereux et al., 1990). However, the relationship between these two metrics is not linear: many bacteria with highly similar 16S genes have hybridization rates much lower than 70% (Stackebrandt and Goebel, 1994), which is the traditional cut-off for delineating species. This trend has been corroborated across diverse prokaryotes (Hauben et al., 1997, 1999; Kang et al., 2007). In addition, a meta-analysis of 16S gene sequencing and DNA hybridization data from 45 bacterial genera further clarified these observations (Keswani and Whitman, 2001). This analysis established that 78% of the variability in hybridization rates could be accounted for by 16S similarity, based on a non-linear model. However, they also identified that a minority of hybridization rates were extremely poorly predicted by 16S similarity (Keswani and Whitman, 2001).

that experimental validation of pathway reconstructions is possible (Francke et al., 2005; Oberhardt et al., 2008). Such validations would be possible if predictions are based on individual members of a microbiome (e.g., metagenome-assembled genomes), but it is less clear what experiments could validate pathways predicted for an overall community. In MGS data pathways are typically inferred as though all gene families were free to interact with each other. In other words, they are inferred as though there was universal cross-feeding. All three approaches described above are intended to be used for such community-wide gene family profiles. However, as mentioned above, this assumption is invalid because clearly not all proteins and metabolites in the microbiome can freely interact (McMahon, 2015). The implications of this assumption being invalid remain unclear, but nonetheless it is an important caveat when interpreting pathway reconstruction data based on community-wide MGS data.

Marker-gene-based metagenome prediction methods

Ideally, analyses of microbial functions are based on MGS data. However, predicted functions based on 16S rRNA gene sequencing data are often analysed instead. [Metagenome prediction](#), predicting complete genomes for each individual amplicon sequence variant 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 important drawbacks. The key problematic assumption is that the marker-gene used for predictions, typically the 16S, is strongly associated with genome content. This broad assumption is correct: genera such as *Lactobacillus* and *Desulfobacter* can be easily distinguished based on the 16S and they are enriched for extremely different functions. Namely, *Lactobacillus* can often perform lactic acid fermentation (Duar et al., 2017) whereas *Desulfobacter* can typically oxidize acetate to CO₂ (Galushko and Kuever, 2019). Such comparisons of characteristic functions between distantly related taxa are uncontroversial. The difficulty arises when approaches attempt to predict entire genome contents for an entire community, including for closely related taxa.

Early work identified substantial genomic variation between closely related taxa, including those with highly similar 16S sequences (see Box 3). These observations agree well with genomic comparisons of strains, which can drastically differ in genome content. For example, across 17 *E. coli* genomes there are 13,000 genes that are variably distributed and only 2,200 core genes (Rasko et al., 2008). This enormous range of genomic variation is not reflected at the 16S level, where *E. coli* strains are typically >99% identical (Suardana, 2014). These genomic differences can translate to enormous variation at higher taxonomic levels as well. For instance, a comparison of the genomes from 11 *Yersinia* species found a range of genome sizes from 3.7 - 4.8 megabases (Chen et al., 2010). A closer comparison of three pathogenic species of *Yersinia* determined that they shared 2,558 protein clusters while 2,603 were variably distributed. These species-level differences are also not proportionally reflected by divergence in *Yersinia* species 16S genes, which are typically >97% identical (Ibrahim et al., 1993). These examples highlight that 16S similarity can be a poor predictor of genomic similarity. This issue is compounded when there are divergent 16S copies within the same genome, although typically these are >99.5% identical (Větrovský and Baldrian, 2013).

Variation in gene content within a taxonomic lineage is a recurrent observation across microbial communities. Variably present genes are often linked to putative niche-specific adaptations (Wilson et al., 2005), such as genes affecting antibiotic resistance (Kallonen et al., 2017), carbohydrate catabolism (Arboleña et al., 2018), and wound healing (Kalan et al., 2019). Based on

these and other observations, the understanding of bacterial genomic content has shifted from that of a static genome to a pan-genome, consisting of core and variable genes (Tettelin et al., 2005). Variably present genes are transmitted between genomes through horizontal gene transfer, which typically occurs between closely related organisms (Popa and Dagan, 2011). However, horizontal gene transfer can also occur between distantly related organisms, such as between different bacterial phyla (Beiko et al., 2005; Kloesges et al., 2011; Martiny et al., 2013).

The high variability between bacterial genomes and extensive horizontal gene transfer highlights the major challenges facing metagenome prediction. Despite these challenges, interest in performing metagenome predictions has continued, supported by several observations. First, although there are important outliers, 16S sequence identity does logarithmically correlate well with the average nucleotide identity between genomes, with an R^2 of 0.79 (Konstantinidis and Tiedje, 2005). Second, 16S sequence similarity does provide some information on the ecological similarity of bacteria (Chaffron et al., 2010). This was demonstrated by the fact that co-occurring environmental bacteria are more likely to have similar 16S sequences. In addition, overall differences in inferred KEGG pathway potential are strongly associated with 16S divergence (Chaffron et al., 2010). Last, within a given environment, such as the human gut, 16S divergence was shown to be particularly predictive of divergence in average gene content (Zaneveld et al., 2010).

Originally, metagenome prediction workflows were based on matching 16S sequences to reference genomes. By taking the best matching genome or averaging across genomes with similar sequences, a predicted genome annotation can be acquired for all 16S sequences (Figure 3). To infer the metagenome profile one must simply multiply the predicted genome annotations for each 16S sequence by the abundance of each 16S sequence in the metagenome. In addition to predicting microbial functions linked to Crohn's disease (Morgan et al., 2012), this approach has also been used to profile diet-related microbial functions across mammals (Muegge et al., 2011) and the functions of invasive bacteria within corals (Barott et al., 2012). Although bioinformatics tools for metagenome prediction are now typically used for performing this task, this 16S-matching approach is still used for custom analyses (Bradley and Pollard, 2020; Verster and Borenstein, 2018).

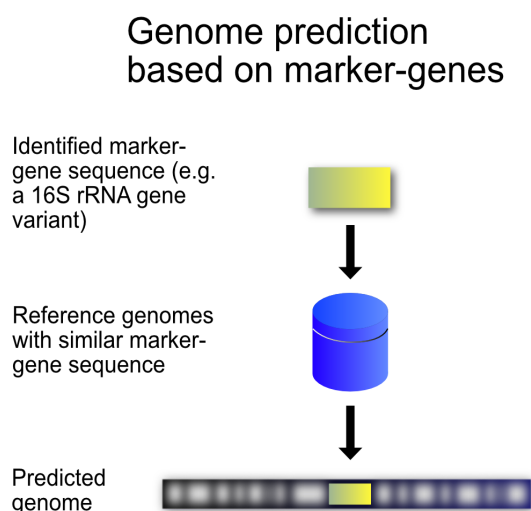


Figure 3: Genome prediction based on marker-gene sequences. This is another method of producing joint taxa and function profiles, which in this case are explicitly linked, similar to assembling genomes. This is also the first step in most metagenome prediction workflows. However, all such genome prediction methods 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.

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

PICRUSt1 introduced the step of normalizing relative abundances by the predicted number of 16S copies within each genome, which is intended to control biases in 16S sequencing due to copy number (Farrelly et al., 1995). Importantly, although 16S copy number correction has become a common step for metagenome prediction (Angly et al., 2014), accurately predicting 16S copy number is particularly challenging. An independent validation of several 16S copy number prediction methods, including PICRUSt1, identified poor agreement of predicted copy numbers against existing reference genomes (Louca, Doebeli, and Parfrey, 2018). In some cases, less than 10% of the variance in actual 16S copy number was explained by these predictions. In addition, these predictions were often only slightly correlated between prediction methods.

Since PICRUSt1 was published a number of similar metagenome prediction tools have been developed. All of these approaches aim to capture the shared phylogenetic signal in the distribution of functions across taxa. These tools include: PanFP (Jun et al., 2015), Piphillin (Iwai et al., 2016; Narayan et al., 2020), PAPRICA (Bowman and Ducklow, 2015), and Tax4Fun2 (Wemheuer et al., 2020).

These metagenome prediction tools have primarily been validated by comparing how well the predicted gene family abundances they output correlate with the abundances of gene families identified in MGS data from the same samples. This approach generally identifies high correlations between the two profiles. For example, predicted KEGG orthologs output by PICRUSt1 based on Human Microbiome Project samples were highly correlated with the matching MGS-identified data (Spearman's $\rho = 0.82$) (Langille et al., 2013). Importantly, a high Spearman correlation is actually expected by chance in these comparisons simply because many genes are common in most environments while others are usually absent or rare. Upon comparing to this expectation the predictions are still significantly better than expected by chance, but only slightly (Douglas et al., 2020). Nonetheless, based on this approach, we found that PICRUSt2 performed marginally better than other tools (Douglas et al., 2020). However, it is noteworthy that Piphillin, which represents a much simpler approach based on a nearest-neighbour approach, performed only slightly worse overall and better in some contexts.

An alternative approach for evaluating these methods is based on the concordance of differential abundance results between actual and predicted metagenomics profiles. When we conducted this analysis while validating PICRUSt2, we found that differential abundances tests on metagenome prediction tools agreed only moderately well with matching tests based on actual MGS data (Douglas et al., 2020). This is a crucial point to appreciate when analyzing metagenome prediction data; even though the overall predicted profiles might correlate with MGS profiles, the results from differential abundance testing might nonetheless be quite different. We also observed high variation across datasets in concordance between MGS and 16S-based predictions. In other words, differential abundance testing on predicted profiles resulted in fair agreement with MGS data on some datasets while disagreeing almost entirely on others. In addition, researchers performing independent work in this area have identified conflicting signals of how well individual metagenome prediction tools perform (Narayan et al., 2020; Sun et al., 2020). These observations might again reflect the high variation across datasets in how well prediction profiles agree with MGS results.

Current state of the integration of taxonomic and functional data types

The above discussion has described the many faces of microbiome data types. Taxonomic and functional microbiome data are typically generated independently, but in some cases can be directly linked (e.g., in metagenome-assembled genomes). Regardless of the exact processing workflow for these data types, we have yet to address one question: how are they integrated?

For independent taxonomic and functional data types this is largely done anecdotally. For example, this is commonly done in regards to the nine genera that are the primary producers of short-chain fatty acids in the human gut (Moya and Ferrer, 2016). Short-chain fatty acid levels have long had an ambiguous link with Crohn's disease (Treem et al., 1994), although they are typically negatively associated with disease activity (Venegas et al., 2019). Due to this association, there has been long-standing interest in identifying microbial taxa that are associated with altered short-chain fatty acid levels. Accordingly, Crohn's disease microbiome studies commonly hypothesize that shifts in the relative abundance of any known short-chain fatty acid-producing taxa likely cause altered short-chain fatty acid levels. For example, *Faecalibacterium prausnitzii* is a well-known commensal short-chain fatty acid-producer in the human gut and is consistently found at lower levels in Crohn's disease patient microbiomes (Wright et al., 2015). Although potential links between lower levels of this species, in addition to other taxa such as *Roseburia* (Laserna-Mendieta et al., 2018), and short-chain fatty acid levels are often discussed, this is rarely formally investigated.

More often, anecdotal links between function and taxa are based on observed associations between significant features. Several such cases have previously been noted as representative examples (Manor and Borenstein, 2017b). For instance, *Propionibacterium acnes* has been identified as strongly correlated with NADH dehydrogenase levels in the skin microbiome (Oh et al., 2014). Consequently, this species was implicated as the likely cause for changes in NADH dehydrogenase levels. Similarly, *Bacteroides thetaiotaomicron* relative abundance has been identified as positively correlated with microbial genes involved with the degradation of complex sugars and starch in the infant gut (Bäckhed et al., 2015). Based on this observation, this species was hypothesized to be the key contributor to increased levels of these degradation genes. Such insights are valuable, but as previously discussed (Manor and Borenstein, 2017b), these anecdotal links alone are not convincing evidence that particular taxa are the primary contributors to functional shifts.

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

In contrast to these examples, there have been calls for improved integration of these microbiome data types, which is rooted in a systems-level biology outlook (Greenblum et al., 2013). "Functional Shifts' Taxonomic Contributors" (FishTaco) is one bioinformatics method developed for this purpose, which quantifies taxonomic contributions to functional shifts (Manor and Borenstein, 2017b). One major application of this approach is to distinguish two explanations for why a function might be at high relative abundance (Figure 4). First, a function might be higher in relative abundance simply because it hitchhiked on the genome of a taxon that bloomed for other reasons. In contrast, an alternative explanation might be that many taxa performing the same function gained a growth advantage and thus grew in relative abundance. FishTaco can also identify functions that have grown in relative abundance simply because microbes that do not encode it are at lower levels.

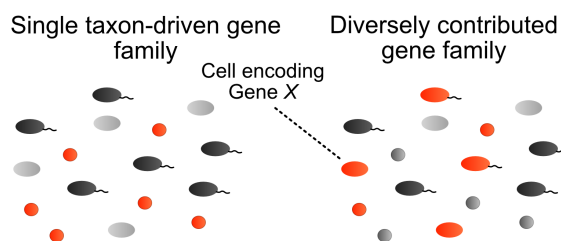


Figure 4: Two explanations for why a gene family might be at higher relative abundance that would be impossible to distinguish without joint taxonomic and functional data. Microbes encoding the gene of interest (Gene X) are indicated in red. This diagram contrasts how a gene family might be blooming 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.

FishTaco works by first identifying significant shifts in functional abundances with a standard differential abundance test, typically a Wilcoxon test. Subsequently, a permutation analysis is undertaken, which consists of randomly shifting the relative abundance of a subset of taxa, while maintaining the rest. A large collection of such permutations is performed, which include permutations of single and multiple taxa in different replicates. Based on this approach an estimate of the relative contribution of each taxon to a functional shift can be calculated (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 compositionally sound analysis and integrate taxonomic information when identifying significant functions.

An alternative method is *phylogenize*, which does address each of these issues (Bradley and Pollard, 2020; Bradley et al., 2018). This approach tests for significant associations between the presence of a taxon 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 across multiple phyla that might be linked with a condition of interest. From another perspective, this focus on the phylum level is also not specific enough; although phylum-function associations are valuable they do not provide information on the relative contributions of lower-level taxa, such as species, to the association. Accordingly, there is room for improvement in both the statistical analysis and interpretation of the *phylogenize* approach.

Despite the availability of approaches for integrating functional and taxonomic data, they have yet to become a mainstay of microbiome analyses. However, it is becoming common to visualize stacked bar-charts of taxonomic contributors to functions of interest (see Figure 5 for examples), which can be created with tools such as BURRITO (McNally et al., 2018). This is typically performed on predicted metagenome output by PICRUSt or alternatively on HUMAnN2 output, although this could be performed with any linked taxa-function data. As discussed above, the HUMAnN2 pipeline includes a step for identifying particular strains in MGS dataset, which allows gene families to be linked to those strains (Franzosa et al., 2018). In some cases this approach enables complete links between taxa and function to be identified. For instance, *F. prausnitzii* was shown to be the obvious principal contributor to glutaryl-CoA biosynthesis in the Human Microbiome Project gut MGS samples (Franzosa et al., 2018). However, more commonly there are numerous taxonomic contributors to a single given function, and it is difficult to interpret which taxa are the key contributors by looking at visualizations alone. Nonetheless, even in the presence of many taxonomic contributors, the HUMAnN2 authors demonstrated that these visualizations can provide information about the diversity of taxa contributing to a function, termed the **contributonal diversity** (Franzosa et al., 2018). This is most often quantified with the Gini-Simpson index, which is the complement of Simpson's evenness (Jost, 2006).

Contributonal diversity has been shown to be a useful approach for delineating housekeeping pathways encoded by many taxa, intermediate pathways, and those rarely encoded, which can correspond to opportunists or keystone species (Figure 5). For instance, *F. prausnitzii* has previously been linked with several human microbiome pathways identified through MGS that have intermediate contributonal diversities (Abu-Ali et al., 2018). When present, this species tended to contribute the majority of all pathways it encoded.

This approach has also been valuable for profiling shifts in the contributions to microbial pathways over time, such as in the infant gut profiled with MGS (Vatanen et al., 2018). In this case, several microbial pathways, such as siderophore biosynthesis, were found to display decreasing contributonal diversity with age. This is an interesting observation because siderophores are costly to produce but are highly beneficial in the human gut. In particular, siderophores can confer a strong benefit to multiple community members, including those that do not produce siderophores, by providing access to iron. Siderophores have previously been presented as microbial functions whose distribution is consistent with the Black Queen Hypothesis (Morris et al., 2012). This hypothesis states that adaptive gene loss may occur for functions that are costly to produce, provided that the function is provided by other community members. This hypothesis was discussed in the context of the infant microbiome as an explanation for why siderophore

contributonal diversity decreases over time (Vatanen et al., 2018): perhaps gene loss confers an adaptive benefit by avoiding the production of a costly metabolite. Although this is an interesting hypothesis, a less controversial interpretation of this result is simply that siderophores became less stably encoded over time in the profiled samples.

Related to this point, two additional metrics have also been developed to summarize the stability of taxonomic contributions to microbial functions (Eng and Borenstein, 2018). More specifically, these metrics are intended to summarize functional robustness across samples, which is the stability in the relative abundance for a given function in response to taxonomic perturbation. This is performed by generating a taxa-response curve that describes the average change in functional relative abundances in response to taxonomic perturbations of different magnitudes. Two metrics are then computed based upon these curves: attenuation and buffering. Attenuation captures how rapidly a function shifts with increasing taxonomic perturbation magnitudes. In contrast, buffering represents how well functional shifts are suppressed at smaller taxonomic perturbation magnitudes.

Applying these metrics to PICRUSt-predicted metagenomes from 16S sequencing of human body sites, validated by a subset of MGS samples, yielded several novel perspectives. First, attenuation and buffering were conserved across body sites for microbial house-keeping pathways but varied for several others. For instance, robustness in the biosynthesis of unsaturated fatty acids varied substantially across body sites. In addition, human gut samples were found to have higher values of both attenuation and buffering than compared to vaginal samples. These trends were shown to be driven by more than simply lower richness in vaginal samples by subsampling to comparable diversity levels across each body-site (Eng and Borenstein, 2018). These observations are consistent with the controversial hypothesis that microbial communities may be under varying selection strengths for functional robustness, depending on the environment (Ley et al., 2006; Naeem et al., 1998).

The development of these metrics for summarizing functional contributions represent an important goal of microbiome research, which is to leverage sequencing data to yield novel biological insights. In contrast, another major goal is to answer a more practical question: how useful is microbiome data for classification and prediction tasks?

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

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

Despite this great interest in applying machine learning to different microbiome data types, there has been little focus on integrating across them. The aforementioned comparison of 16S-based taxa and predicted functions is one exception where a hybrid classification model of both

data types was created (Ning and Beiko, 2015). In this case, there was a small increase in classification performance for distinguishing nine human oral sub-locations. The original operational taxonomic unit and KEGG ortholog-based models yielded accuracies of 76.2% and 76.1%, respectively, while the hybrid model resulted in an accuracy of 77.7% (Ning and Beiko, 2015). This result indicates that predicted functions may provide some additional information in combination with taxonomic data, but the consistency and biological significance of this small effect remains unclear. Further investigation into the integration of these data types within a machine learning context is needed to ensure that the highest-quality models possible are constructed.

Outlook

Herein we have described the unique characteristics of microbiome DNA data types and many of the approaches that have been proposed for their analysis. Throughout we have emphasized two ideas. First, increased integration of taxonomic and functional microbiome data types is needed. And second, there is often high variation in the results between microbiome data analysis pipelines.

Regarding the first point, we believe that several of the tools described above, such as Fish-Taco and *phylogenize*, largely solve the issue of how to jointly investigate taxa and functions. Increased usage and development of these and other related tools would greatly help with the interpretability of microbiome data.

One area where further development is particularly needed is in the context of classification models, where little work has been conducted to systematically link taxa and functions appropriately. One exception was a classification approach based on gene families that identified predictive genes and then subsequently identified metagenome-assembled genomes within a given dataset enriched for these genes (Rahman et al., 2018). However, this approach still relied on follow-up analyses rather than integrating the data types. Instead, an improved approach could be based on explicitly leveraging the hierarchical nature of microbiome data types. This is because functional and taxonomic data types independently form clear hierarchical structures (e.g., Pathway - Gene and Phylum - Class - Order). The connection between taxa and gene families and pathways is more complex, but nonetheless, links between groups of strains or amplicon sequence variants and microbial functions can be defined. A modified machine learning framework that explicitly accounted for these relationships could result in more interpretable outputs.

Regardless of the specific tool, microbiome researchers should move towards more integration of taxonomic and functional data. It is odd to distinguish between functional and taxonomic datatypes in the first place: they are inextricably linked after all. The term "metagenome" itself is in some ways unfortunate as it implies that the genetic information for all organisms in a community can be simultaneously analyzed in a coherent way, without partitioning genes into genomes. This may be valid for high-level pathways but for generating hypotheses regarding specific gene families it is too often misleading. This perspective is becoming more common, as the availability of metagenome-assembled genomes increases (Frioux et al., 2020).

The other common thread throughout this manuscript has been that technical variation in microbiome data analyses means that making robust biological inferences, especially regarding specific microbial features, is challenging. Indeed, the lack of standardization in microbiome data analysis has previously been strongly criticized. An assessment of numerous papers attempting to define standard pipelines concluded that there was disturbingly little consensus (Pollock et al., 2018). This is true for many steps related to the processing, sequencing, and analysis of microbiome data (McLaren et al., 2019; Sinha et al., 2017). For instance, there have been contradictory results regarding the efficacy of different extraction protocols (Greathouse et al., 2019; Salonen et al., 2010). In particular, underrepresentation of Gram-positives has been observed (Maukonen et al., 2012), which may be partially resolved by using bead-beating extraction protocols (Guo and Zhang, 2013). Common extraction protocols also often result in high rates of DNA fragmentation, which makes the extracted DNA less appropriate for long-read sequencing technologies. Updated extraction protocols based on robust enzymatic lysis have been developed to address this problem (Maghini et al., 2020). There is also substantial technical variation related to

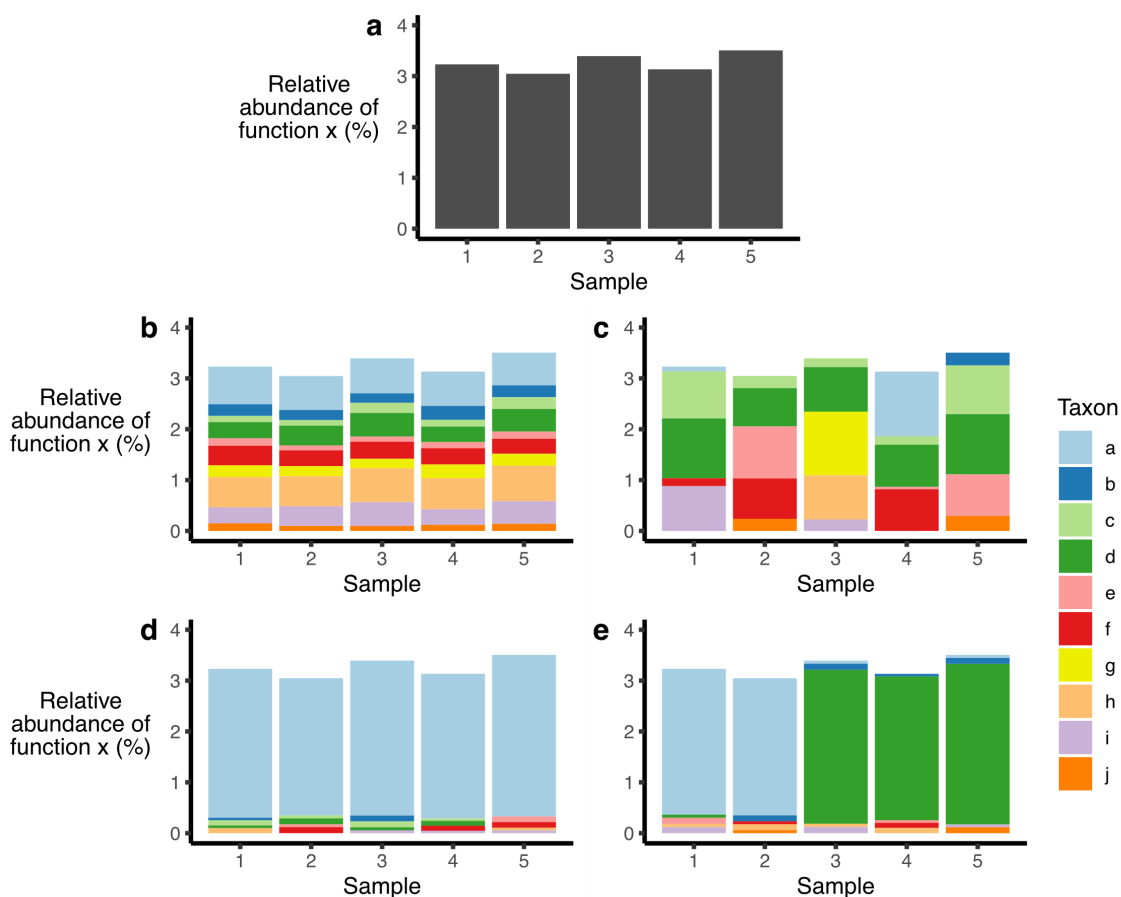


Figure 5: Contrasting extremes of taxa contributing to a function across samples. A simple example of relative abundances for a given microbial function (e.g., a pathway) across five samples, which is at similar levels across all samples. It is common to analyze microbial functions without integrating taxonomic information. However, including this taxonomic information can help distinguish many different biological scenarios. Panels b–e represent four extreme scenarios that could account for the relative abundances shown in panel a. These examples also highlight the value of using stacked bar charts and are closely based on the examples presented by Franzosa et al., 2018. The examples: (b) stable contribution of the function by the same diverse taxa; (c) contribution of the taxa by different taxa; (d) stable contribution primarily by a single taxon; and (e) contribution primarily by a single taxon, which can differ across samples.

bioinformatics choices, which represent the final steps of a microbiome project. For example, as discussed above, the bioinformatics choices made when performing differential abundance testing on microbiome data can have severe impacts on any interpretations (Hawinkel et al., 2019; Thorsen et al., 2016).

We have encountered similar issues with our work, most strikingly when investigating pediatric Crohn's disease patients' microbiome profiles (Douglas et al., 2018). An important characteristic of these data was that 98% of the sequenced reads mapped to the human genome. This characteristic made taxonomic profiling of these data especially prone to false positives. In particular, an initial draft of our manuscript was based on profiles that included large proportions of viral-identified DNA and matches to certain eukaryotic parasites. We were initially excited about these observations, because the abundances of these non-prokaryotic taxa were discriminative for classifying patient disease state and treatment response. However, the exact taxa identified were peculiar: they were predominately represented by a range of plant-associated viruses and the eukaryotic genus *Plasmodium*, which is best known as including the causative agent for

malaria, *Plasmodium falciparum*. Upon closer investigation it became clear that this signal was driven entirely by a difference in how reads were mapped to lineage-specific marker-genes with MetaPhlan2. Altering the parameter choice from local to global mapping entirely removed these taxa. This relatively small difference in parameter choice appeared to only affect our data and not more typical microbiome datasets, which we believe was due to the high proportion of human DNA in our data. Although this error was moderately embarrassing, it was more importantly an example of how easily a single parameter setting can result in starkly different biological interpretations. In this case the difference was driven by an option used for a single bioinformatics tool.

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

Similarly lower alpha diversity in stool microbiomes has been frequently linked with disease states (Mosca et al., 2016). These observations are intuitively reasonable as reduced alpha diversity could enable pathogens to bloom (Vincent et al., 2013) or represent differences in resource availability (Turnbaugh et al., 2009). However a re-analysis of data from 28 studies representing ten diseases was unable to identify evidence for links between alpha diversity and disease states (Duvall et al., 2017). The exceptions were diarrheal diseases and inflammatory bowel diseases.

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

Until a clear consensus is reached it is the responsibility of microbiome researchers to make the caveats and challenges facing this area clear to readers and newcomers to the field. This is crucial given the widespread interest in studying microbiomes through DNA sequencing; the number of microbiome sequencing-related publications continues to rapidly grow. This is in tandem with funding for these projects, which has steadily increased in the USA from at least 2007 to 2016 (NIH, 2019). According to the US National Health Institute, there was US\$766 million dollars invested in microbiome research in 2019, which was the 63rd most highly funded health-related research category out of 291. Although comparing across research categories of varying granularity is difficult, it is noteworthy that microbiome research was more highly funded than both breast cancer and Alzheimer's disease research. Importantly, an increased interest in microbiome research is warranted: recent technological developments are enabling improved investigations into microbial biology. However, as the monetary investment and research hours dedicated to microbiome research grows, it is crucial that scientists ensure the best use of these resources. Open discussions on the many contentious aspects of microbiome data analysis would help with this issue. Indeed, such clarifications by leaders in the microbiome field are starting become more common (Allaband et al., 2019). However, although these contributions are valuable, they do not adequately address the problem. In particular, instead of mentioning these issues in passing, inconsistencies between bioinformatics workflows should be emphasized more clearly for the benefit of the uninitiated.

Another practical improvement would be to normalize, and potentially require, explicit summaries of the effects of technical variation on any biological interpretations reported in microbiome studies. This is impossible to capture entirely, but it could be done by comparing how key results change depending on a subset of representative bioinformatics choices. For instance, researchers could compare how insights change depending on the combinations of denoising tools and differential abundance methods that they have applied when analyzing 16S data. Although these changes would result in increased workloads when conducting analyses and when communicating results, they would help ensure that any major biological findings are at least robust to a representative set of bioinformatics choices.

Regardless of which approach is taken to address these issues, the most important point is that action is needed on this front. The variation between bioinformatics methods is undeniable and unfortunately reflects a reproducibility crisis facing microbiome data analysis.

Acknowledgements

We would like to thank the following individuals for feedback on sections of this manuscript: Dr. Robert Beiko, Dr. Zhenyu Cheng, Dr. André Comeau, Casey Jones, Jacob Nearing, Dr. Laura Parfrey, Dr. Andrew Stadnyk, Chris Tang, and Dr. Robyn Wright. A previous version of this article was peer-reviewed and recommended by Peer Community In Genomics (<https://doi.org/10.24072/pci.genomics.100008>). We would like to thank the recommender and the three reviewers who provided helpful feedback on our preprint while it was under review.

Conflict of interest disclosure

The authors of this article declare that they have no financial conflict of interest with the content of this article.

References

- Abellan-Schneyder I et al. (2021). *Primer, Pipelines, Parameters: Issues in 16S rRNA Gene Sequencing*. *mSphere* **6**, e01202–20. <https://doi.org/10.1128/mSphere.01202-20>.
- Abu-Ali GS et al. (2018). *Metatranscriptome of human faecal microbial communities in a cohort of adult men*. *Nature Microbiology* **3**, 356–366. <https://doi.org/10.1038/s41564-017-0084-4>.
- Abubucker S et al. (2012). *Metabolic reconstruction for metagenomic data and its application to the human microbiome*. *PLoS Computational Biology* **8**, e1002358. <https://doi.org/10.1371/journal.pcbi.1002358>.
- Aitchison J (1982). *The Statistical Analysis of Compositional Data*. *Journal of the Royal Statistical Society. Series B (Methodological)* **44**, 139–177. <https://doi.org/10.1111/j.2517-6161.1982.tb01195.x>.
- Allaband C et al. (2019). *Microbiome 101: Studying, Analyzing, and Interpreting Gut Microbiome Data for Clinicians*. *Clinical Gastroenterology and Hepatology* **17**, 218–230. <https://doi.org/10.1016/j.cgh.2018.09.017>.
- Ambler RP et al. (1979). *Cytochrome C2 sequence variation among the recognised species of purple nonsulphur photosynthetic bacteria*. <https://doi.org/10.1038/278659a0>.
- Amir A et al. (2017). *Deblur Rapidly Resolves Single- Nucleotide Community Sequence Patterns*. *mSystems* **2**, e00191–16. <https://doi.org/10.1128/mSystems.00191-16>.
- Angly FE et al. (2014). *CopyRighter: A rapid tool for improving the accuracy of microbial community profiles through lineage-specific gene copy number correction*. *Microbiome* **2**, 11. <https://doi.org/10.1186/2049-2618-2-11>.
- Apweiler R et al. (2004). *UniProt: the Universal Protein knowledgebase*. *Nucleic acids research* **32**, D115–119. <https://doi.org/10.1093/nar/gkh131>.
- Arboleya S et al. (2018). *Gene-trait matching across the Bifidobacterium longum pan-genome reveals considerable diversity in carbohydrate catabolism among human infant strains*. *BMC Genomics* **19**, 33. <https://doi.org/10.1186/s12864-017-4388-9>.

- Armour CR et al. (2019). A Metagenomic Meta-analysis Reveals Functional Signatures of Health and Disease in the Human Gut Microbiome. *mSystems* **4**, e00332–18. <https://doi.org/10.1128/msystems.00332-18>.
- Ayling M et al. (2020). New approaches for metagenome assembly with short reads. *Briefings in Bioinformatics* **21**, 584–594. <https://doi.org/10.1093/bib/bbz020>.
- Bäckhed F et al. (2015). Dynamics and stabilization of the human gut microbiome during the first year of life. *Cell Host and Microbe* **17**, 690–703. <https://doi.org/10.1016/j.chom.2015.04.004>.
- Bahram M et al. (2019). Newly designed 16S rRNA metabarcoding primers amplify diverse and novel archaeal taxa from the environment. *Environmental Microbiology Reports* **11**, 487–494. <https://doi.org/10.1111/1758-2229.12684>.
- Barott KL et al. (2012). Microbial to reef scale interactions between the reef-building coral *Montastraea annularis* and benthic algae. *Proceedings of the Royal Society B: Biological Sciences* **279**, 1655–1664. <https://doi.org/10.1098/rspb.2011.2155>.
- Beiko RG et al. (2005). Highways of gene sharing in prokaryotes. *PNAS USA* **102**, 14332–14337. <https://doi.org/10.1073/pnas.0504068102>.
- Bellemain E et al. (2010). ITS as an environmental DNA barcode for fungi: an in silico approach reveals potential PCR biases. *BMC Microbiology* **10**, 189. <https://doi.org/10.1186/1471-2180-10-189>.
- Belton JM et al. (2012). Hi-C: A comprehensive technique to capture the conformation of genomes. *Methods* **58**, 268–276. <https://doi.org/10.1016/j.ymeth.2012.05.001>.
- Berg G et al. (2020). Microbiome definition re-visited: old concepts and new challenges. *Microbiome* **8**, 103. <https://doi.org/10.1186/s40168-020-00875-0>.
- Bishara A et al. (2018). High-quality genome sequences of uncultured microbes by assembly of read clouds. *Nature Biotechnology*. <https://doi.org/10.1038/nbt.4266>.
- Bowers RM et al. (2017). Minimum information about a single amplified genome (MISAG) and a metagenome-assembled genome (MIMAG) of bacteria and archaea. *Nature Biotechnology* **35**, 725–731. <https://doi.org/10.1038/nbt.3893>.
- Bowman JS, Ducklow HW (2015). Microbial communities can be described by metabolic structure: A general framework and application to a seasonally variable, depth-stratified microbial community from the coastal West Antarctic Peninsula. *PLoS ONE* **10**, e0135868. <https://doi.org/10.1371/journal.pone.0135868>.
- Bradley PH, Pollard KS (2020). Phylogenize: Correcting for phylogeny reveals genes associated with microbial distributions. *Bioinformatics* **36**, 1289–1290. <https://doi.org/10.1093/bioinformatics/btz722>.
- Bradley PH et al. (2018). Phylogeny-corrected identification of microbial gene families relevant to human gut colonization. *PLoS Computational Biology* **14**, e1006242. <https://doi.org/10.1371/journal.pcbi.1006242>.
- Breiman L (2001). Random Forests. *Machine Learning* **45**, 5–32. <https://doi.org/10.1023/A:1010933404324>.
- Breitwieser FP et al. (2019). A review of methods and databases for metagenomic classification and assembly. *Briefings in Bioinformatics* **20**, 1125–1139. <https://doi.org/10.1093/bib/bbx120>.
- Brenner DJ (1973). Deoxyribonucleic acid reassociation in the taxonomy of enteric bacteria. *International Journal of Systematic Bacteriology* **23**, 298–307. <https://doi.org/10.1099/00207713-23-4-298>.
- Buchfink B et al. (2015). Fast and Sensitive Protein Alignment using DIAMOND. *Nature Methods* **12**, 59–60. <https://doi.org/10.1038/nmeth.3176>.
- Bukin YS et al. (2019). The effect of 16s rRNA region choice on bacterial community metabarcoding results. *Scientific Data* **6**, 190007. <https://doi.org/10.1038/sdata.2019.7>.
- Burke C et al. (2011). Bacterial community assembly based on functional genes rather than species. *Proceedings of the National Academy of Sciences of the USA* **108**, 14288–14293. <https://doi.org/10.1073/pnas.1101591108/-/DCSupplemental.www.pnas.org/cgi/doi/10.1073/pnas.1101591108>.

- Calgaro M et al. (2020). Assessment of single cell RNA-seq statistical methods on microbiome data. *Genome Biology*, 191. <https://doi.org/10.1101/2020.01.15.907964>.
- Callahan BJ, Sankaran K, et al. (2016). Bioconductor workflow for microbiome data analysis: From raw reads to community analyses. *F1000 Research* 5, 1492. <https://doi.org/10.12688/F1000RESEARCH.8986.1>.
- Callahan BJ, McMurdie PJ, et al. (2016). DADA2: High resolution sample inference from amplicon data. *Nature Methods* 13, 581–583. <https://doi.org/10.1101/024034>.
- Callahan BJ et al. (2019). High-throughput amplicon sequencing of the full-length 16S rRNA gene with single-nucleotide resolution. *Nucleic acids research* 47, e103. <https://doi.org/10.1093/nar/gkz569>.
- Carini P et al. (2016). Relic DNA is abundant in soil and obscures estimates of soil microbial diversity. *Nature Microbiology* 2, 16242. <https://doi.org/10.1038/nmicrobiol.2016.242>.
- Case RJ et al. (2007). Use of 16S rRNA and rpoB genes as molecular markers for microbial ecology studies. *Applied and Environmental Microbiology* 73, 278–288. <https://doi.org/10.1128/AEM.01177-06>.
- Caspi R et al. (2013). The MetaCyc database of metabolic pathways and enzymes and the BioCyc collection of pathway/genome databases. *Nucleic Acids Research* 42, D459–D471. <https://doi.org/10.1093/nar/gkv1164>.
- Castaner O et al. (2018). The gut microbiome profile in obesity: A systematic review. *International Journal of Endocrinology* 2018, 4095789. <https://doi.org/10.1155/2018/4095789>.
- Chaffron S et al. (2010). A global network of coexisting microbes from environmental and whole-genome sequence data. *Genome Research* 20, 947–959. <https://doi.org/10.1101/gr.104521.109>.
- Chen PE et al. (2010). Genomic characterization of the Yersinia genus. *Genome Biology* 11, R1. <https://doi.org/10.1186/gb-2010-11-1-r1>.
- Chen Z et al. (2019). Impact of Preservation Method and 16S rRNA Hypervariable Region on Gut Microbiota Profiling. *mSystems* 4, e00271–18. <https://doi.org/10.1128/msystems.00271-18>.
- Comeau AM et al. (2017). Microbiome Helper: a Custom and Streamlined Workflow for Microbiome Research. *mSystems* 2, e00127–16. <https://doi.org/10.1128/mSystems.00127-16>.
- Comeau AM et al. (2019). Nitrate Consumers in Arctic Marine Eukaryotic Communities: Comparative Diversities of 18S rRNA, 18S rRNA Genes, and Nitrate Reductase Genes. *Applied and environmental microbiology* 85, e00247–19. <https://doi.org/10.1128/AEM.00247-19>.
- Cruz GNF et al. (2021). Equivolumetric protocol generates library sizes proportional to total microbial load in next-generation sequencing. *Frontiers in Microbiology* 12, 638231. <https://doi.org/10.1101/2020.02.03.932301>.
- Darling AE et al. (2014). PhyloSift: Phylogenetic analysis of genomes and metagenomes. *PeerJ* 2014, 243. <https://doi.org/10.7717/peerj.243>.
- Devereux R et al. (1990). Diversity and origin of Desulfovibrio species: Phylogenetic definition of a family. *Journal of Bacteriology* 172, 3609–3619. <https://doi.org/10.1128/jb.172.7.3609-3619.1990>.
- Doolittle WF, Booth A (2017). It's the song, not the singer: an exploration of holobiosis and evolutionary theory. *Biology and Philosophy* 32, 5–24. <https://doi.org/10.1007/s10539-016-9542-2>.
- Douglas GM et al. (2020). PICRUSt2 for prediction of metagenome functions. *Nature Biotechnology* 38, 685–688. <https://doi.org/10.1038/s41587-020-0548-6>.
- Douglas GM, Langille MG (2019). Current and promising approaches to identify horizontal gene transfer events in metagenomes. *Genome Biology and Evolution* 11, 2750–2766. <https://doi.org/10.1093/gbe/evz184>.
- Douglas GM et al. (2018). Multi-omics differentially classify disease state and treatment outcome in pediatric Crohn's disease. *Microbiome* 6, 13. <https://doi.org/10.1186/s40168-018-0398-3>.
- Drouin G et al. (2011). The genetics of vitamin C loss in vertebrates. *Current genomics* 12, 371–8. <https://doi.org/10.2174/138920211796429736>.

- Duar RM et al. (2017). *Lifestyles in transition: evolution and natural history of the genus Lactobacillus*. *FEMS microbiology reviews* **41**, S27–S48. <https://doi.org/10.1093/femsre/fux030>.
- Duvallet C et al. (2017). *Meta-analysis of gut microbiome studies identifies disease-specific and shared responses*. *Nature Communications* **8**, 1784. <https://doi.org/10.1038/s41467-017-01973-8>.
- Edgar RC (2016). *UNOISE2: improved error-correction for Illumina 16S and ITS amplicon sequencing*. *bioRxiv*. <https://doi.org/10.1101/081257>.
- Eng A, Borenstein E (2018). *Taxa-function robustness in microbial communities*. *Microbiome* **6**, 45. <https://doi.org/10.1186/s40168-018-0425-4>.
- Farrelly V et al. (1995). *Effect of genome size and rrn gene copy number on PCR amplification of 16S rRNA genes from a mixture of bacterial species*. *Applied and Environmental Microbiology* **61**, 2798–2801. <https://doi.org/10.1128/aem.61.7.2798-2801.1995>.
- Faust K et al. (2012). *Microbial co-occurrence relationships in the Human Microbiome*. *PLoS Computational Biology* **8**. <https://doi.org/10.1371/journal.pcbi.1002606>.
- Fernandes AD et al. (2014). *Unifying the analysis of high-throughput sequencing datasets: characterizing RNA-seq, 16S rRNA gene sequencing and selective growth experiments by compositional data analysis*. *Microbiome* **2**, 15. <https://doi.org/10.1186/2049-2618-2-15>.
- Fernandes AD et al. (2013). *ANOVA-Like Differential Expression (ALDEx) Analysis for Mixed Population RNA-Seq*. *PLoS ONE* **8**, e67019. <https://doi.org/10.1371/journal.pone.0067019>.
- Finn RD et al. (2014). *Pfam: The protein families database*. *Nucleic Acids Research* **42**, D222–D230. <https://doi.org/10.1093/nar/gkt1223>.
- Fitch WM, Margoliash E (1967). *Construction of phylogenetic trees*. *Science* **155**, 279–284. <https://doi.org/10.1126/science.155.3760.279>.
- Fox GE et al. (1992). *How close is close: 16S rRNA sequence identity may not be sufficient to guarantee species identity*. *International Journal of Systematic Bacteriology* **42**, 166–170. <https://doi.org/10.1099/00207713-42-1-166>.
- Francke C et al. (2005). *Reconstructing the metabolic network of a bacterium from its genome*. *Trends in Microbiology* **13**, 550–558. <https://doi.org/10.1016/j.tim.2005.09.001>.
- Franzosa EA et al. (2018). *Species-level functional profiling of metagenomes and metatranscriptomes*. *Nature Methods* **15**, 962–968. <https://doi.org/10.1038/s41592-018-0176-y>.
- Friedman J, Alm EJ (2012). *Inferring Correlation Networks from Genomic Survey Data*. *PLoS Computational Biology* **8**, e1002687. <https://doi.org/10.1371/journal.pcbi.1002687>.
- Frioux C et al. (2020). *From bag-of-genes to bag-of-genomes: metabolic modelling of communities in the era of metagenome-assembled genomes*. *Computational and Structural Biotechnology Journal* **18**, 1722–1734. <https://doi.org/10.1016/j.csbj.2020.06.028>.
- Galperin MY et al. (2015). *Expanded microbial genome coverage and improved protein family annotation in the COG database*. *Nucleic Acids Research* **43**, D261–D269. <https://doi.org/10.1093/nar/gku1223>.
- Galperin MY et al. (2019). *Microbial genome analysis: The COG approach*. *Briefings in Bioinformatics* **20**, 1063–1070. <https://doi.org/10.1093/bib/bbx117>.
- Galushko A, Kuever J (2019). *Desulfobacter*. John Wiley & Sons, Inc., in association with Bergey's Manual Trust. <https://doi.org/10.1002/9781118960608.gbm01011.pub2>.
- Gevers D et al. (2014). *The treatment-naive microbiome in new-onset Crohn's disease*. *Cell Host and Microbe* **15**, 382–392. <https://doi.org/10.1016/j.chom.2014.02.005>.
- Gloor GB et al. (2016). *It's all relative: analyzing microbiome data as compositions*. **26**. <https://doi.org/10.1016/j.annepidem.2016.03.003>.
- Gloor GB et al. (2017). *Microbiome datasets are compositional: And this is not optional*. *Frontiers in Microbiology* **8**, 2224. <https://doi.org/10.3389/fmicb.2017.02224>.
- Goodrich JK et al. (2014). *Conducting a microbiome study*. *Cell* **158**, 250–262. <https://doi.org/10.1016/j.cell.2014.06.037>.
- Graspeuntner S et al. (2018). *Selection of validated hypervariable regions is crucial in 16S-based microbiota studies of the female genital tract*. *Scientific Reports* **8**, 9678. <https://doi.org/10.1038/s41598-018-27757-8>.

- Greathouse KL et al. (2019). DNA extraction for human microbiome studies: The issue of standardization. *Genome Biology* **20**, 212. <https://doi.org/10.1186/s13059-019-1843-8>.
- Greenblum S et al. (2013). Towards a predictive systems-level model of the human microbiome: Progress, challenges, and opportunities. *Current Opinion in Biotechnology* **24**, 810–820. <https://doi.org/10.1016/j.copbio.2013.04.001>.
- Gruber-Vodicka HR et al. (2020). phyloFlash: Rapid Small-Subunit rRNA Profiling and Targeted Assembly from Metagenomes. *mSystems* **5**, e00920–20. <https://doi.org/10.1128/msystems.00920-20>.
- Guo F, Zhang T (2013). Biases during DNA extraction of activated sludge samples revealed by high throughput sequencing. *Applied Microbiology and Biotechnology* **97**, 4607–4616. <https://doi.org/10.1007/s00253-012-4244-4>.
- Haft DH et al. (2003). The TIGRFAMs database of protein families. *Nucleic Acids Research* **31**, 371–373. <https://doi.org/10.1093/nar/gkg128>.
- Hallam SJ et al. (2003). Identification of Methyl Coenzyme M Reductase A (mcrA) Genes Associated with Methane-Oxidizing Archaea. *Applied and Environmental Microbiology* **69**, 5483–5491. <https://doi.org/10.1128/AEM.69.9.5483-5491.2003>.
- Hao X, Chen T (2012). OTU Analysis Using Metagenomic Shotgun Sequencing Data. *PloS one* **7**, e49785. <https://doi.org/10.1371/journal.pone.0049785>.
- Hastie AR et al. (2013). Rapid Genome Mapping in Nanochannel Arrays for Highly Complete and Accurate De Novo Sequence Assembly of the Complex *Aegilops tauschii* Genome. *PLoS ONE* **8**, e55864. <https://doi.org/10.1371/journal.pone.0055864>.
- Hauben L et al. (1997). Comparison of 16S Ribosomal DNA Sequences of All *Xanthomonas* Species. *International Journal of Systematic Bacteriology* **47**, 328–335. <https://doi.org/10.1099/00207713-47-2-328>.
- Hauben L et al. (1999). Genomic diversity of the genus *Stenotrophomonas*. *International Journal of Systematic Bacteriology* **49**, 1749–1760. <https://doi.org/10.1099/00207713-49-4-1749>.
- Hawinkel S et al. (2019). A broken promise: Microbiome differential abundance methods do not control the false discovery rate. *Briefings in Bioinformatics* **20**, 1–12. <https://doi.org/10.1093/bib/bbx104>.
- Hill C (2020). You have the microbiome you deserve. *Gut Microbiome* **1**, e3. <https://doi.org/10.1017/gmb.2020.3>.
- Hillmann B et al. (2018). Evaluating the Information Content of Shallow Shotgun Metagenomics. *mSystems* **3**, e00069–18. <https://doi.org/10.1128/mSystems.00069-18>.
- HMP-consortium (2013). Structure, Function and Diversity of the Healthy Human Microbiome. *Nature* **486**, 207–214. <https://doi.org/10.1038/nature11234>.
- Hug LA, Edwards EA (2013). Diversity of reductive dehalogenase genes from environmental samples and enrichment cultures identified with degenerate primer PCR screens. *Frontiers in Microbiology* **4**, 341. <https://doi.org/10.3389/fmicb.2013.00341>.
- Huson DH et al. (2007). MEGAN analysis of metagenomic data. *Genome Research* **17**, 377–386. <https://doi.org/10.1101/gr.5969107>.
- Ibrahim A et al. (1993). The phylogeny of the genus *Yersinia* based on 16S rDNA sequences. *FEMS Microbiology Letters* **114**, 173–177. <https://doi.org/10.1111/j.1574-6968.1993.tb06569.x>.
- Inkpen AI et al. (2017). The Coupling of Taxonomy and Function in Microbiomes. *Biology and Philosophy* **32**, 1225–1243. <https://doi.org/10.1007/s10539-017-9602-2>.
- Iwai S et al. (2016). Piphillin: Improved prediction of metagenomic content by direct inference from human microbiomes. *PLoS ONE* **11**, e0166104. <https://doi.org/10.1371/journal.pone.0166104>.
- Jackson DA (1997). Compositional data in community ecology: The paradigm or peril of proportions? *Ecology* **78**, 929–940. [https://doi.org/10.1890/0012-9658\(1997\)078\[0929:CDICET\]2.0.CO;2](https://doi.org/10.1890/0012-9658(1997)078[0929:CDICET]2.0.CO;2).
- Janda JM, Abbott SL (2007). 16S rRNA gene sequencing for bacterial identification in the diagnostic laboratory: Pluses, perils, and pitfalls. *Journal of Clinical Microbiology* **45**, 2761–2764. <https://doi.org/10.1128/JCM.01228-07>.

- Jayaprakash TP et al. (2012). Resolution and characterization of distinct *cpn60*-based subgroups of *Gardnerella vaginalis* in the vaginal microbiota. *PLoS ONE* **7**, e43009. <https://doi.org/10.1371/journal.pone.0043009>.
- Jensen LJ et al. (2008). eggNOG: Automated construction and annotation of orthologous groups of genes. *Nucleic Acids Research* **36**, D250–D254. <https://doi.org/10.1093/nar/gkm796>.
- Johnson JS et al. (2019). Evaluation of 16S rRNA gene sequencing for species and strain-level microbiome analysis. *Nature Communications* **10**, 5029. <https://doi.org/10.1038/s41467-019-13036-1>.
- Jones MB et al. (2015). Library preparation methodology can influence genomic and functional predictions in human microbiome research. *Proceedings of the National Academy of Sciences of the United States of America* **112**, 14024–14029. <https://doi.org/10.1073/pnas.1519288112>.
- Jost L (2006). Entropy and diversity. *Oikos* **113**, 363–375. <https://doi.org/10.1111/j.2006.0030-1299.14714.x>.
- Jun SR et al. (2015). PanFP: Pangenome-based functional profiles for microbial communities. *BMC Research Notes* **8**, 479. <https://doi.org/10.1186/s13104-015-1462-8>.
- Kalan LR et al. (2019). Strain- and Species-Level Variation in the Microbiome of Diabetic Wounds Is Associated with Clinical Outcomes and Therapeutic Efficacy. *Cell Host and Microbe* **25**, 641–655. <https://doi.org/10.1016/j.chom.2019.03.006>.
- Kallonen T et al. (2017). Systematic longitudinal survey of invasive *Escherichia coli* in England demonstrates a stable population structure only transiently disturbed by the emergence of ST131. *Genome Research* **27**, 1437–1449. <https://doi.org/10.1101/gr.216606.116>.
- Kanehisa M et al. (2016). KEGG as a reference resource for gene and protein annotation. *Nucleic Acids Research* **44**, D457–D462. <https://doi.org/10.1093/nar/gkv1070>.
- Kang CH et al. (2007). Relationship between genome similarity and DNA-DNA hybridization among closely related bacteria. *Journal of Microbiology and Biotechnology* **17**, 945–951.
- Keinan A et al. (2004). Fair attribution of functional contribution in artificial and biological networks. *Neural Computation* **16**, 1887–1915. <https://doi.org/10.1162/0899766041336387>.
- Keswani J, Whitman WB (2001). Relationship of 16S rRNA sequence similarity to DNA hybridization in prokaryotes. *International Journal of Systematic and Evolutionary Microbiology* **51**, 667–678. <https://doi.org/10.1099/00207713-51-2-667>.
- Kim D et al. (2016). Centrifuge: rapid and accurate classification of metagenomic sequences. *Genome Research* **26**, 1721–1729. <https://doi.org/10.1101/054965>.
- Kloesges T et al. (2011). Networks of gene sharing among 329 proteobacterial genomes reveal differences in lateral gene transfer frequency at different phylogenetic depths. *Molecular Biology and Evolution* **28**, 1057–1074. <https://doi.org/10.1093/molbev/msq297>.
- Knight R et al. (2018). Best practices for analysing microbiomes. *Nature Reviews Microbiology* **16**, 410–422. <https://doi.org/10.1038/s41579-018-0029-9>.
- Knights D et al. (2011). Human-associated microbial signatures: Examining their predictive value. *Cell Host and Microbe* **10**, 292–296. <https://doi.org/10.1016/j.chom.2011.09.003>.
- Konstantinidis KT, Tiedje JM (2005). Genomic insights that advance the species definition for prokaryotes. *Proceedings of the National Academy of Sciences USA* **102**, 2567–2572. [https://doi.org/10.1016/S0040-4020\(01\)97190-X](https://doi.org/10.1016/S0040-4020(01)97190-X).
- Koonin EV, Galperin MY (2003). *Sequence - Evolution - Function: Computational Approaches in Comparative Genomics*. Boston: Kluwer Academic.
- Kurtz ZD et al. (2015). Sparse and Compositionally Robust Inference of Microbial Ecological Networks. *PLoS Computational Biology* **11**, e1004226. <https://doi.org/10.1371/journal.pcbi.1004226>.
- Langille MGI et al. (2013). Predictive functional profiling of microbial communities using 16S rRNA marker gene sequences. *Nature Biotechnology* **31**, 814–821. <https://doi.org/10.1038/nbt.2676>.
- Laserna-Mendieta EJ et al. (2018). Determinants of reduced genetic capacity for butyrate synthesis by the gut microbiome in Crohn's disease and ulcerative colitis. *Journal of Crohn's and Colitis* **12**, 204–216. <https://doi.org/10.1093/ecco-jcc/jjx137>.

- Lau JT et al. (2016). Capturing the diversity of the human gut microbiota through culture-enriched molecular profiling. *Genome Medicine* **8**, 72. <https://doi.org/10.1186/s13073-016-0327-7>.
- Ley RE et al. (2006). Ecological and evolutionary forces shaping microbial diversity in the human intestine. *Cell* **124**, 837–848. <https://doi.org/10.1016/j.cell.2006.02.017>.
- Li D et al. (2015). MEGAHIT: An ultra-fast single-node solution for large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics* **31**, 1674–1676. <https://doi.org/10.1093/bioinformatics/btv033>.
- Links MG et al. (2012). The Chaperonin-60 Universal Target Is a Barcode for Bacteria That Enables De Novo Assembly of Metagenomic Sequence Data. *PLoS ONE* **7**, e49755. <https://doi.org/10.1371/journal.pone.0049755>.
- Liu J et al. (2015). Comparison of ITS and 18S rDNA for estimating fungal diversity using PCR–DGGE. *World Journal of Microbiology and Biotechnology* **31**, 1387–1395. <https://doi.org/10.1007/s11274-015-1890-6>.
- Lloyd-Price J et al. (2017a). Erratum: Strains, functions and dynamics in the expanded Human Microbiome Project. *Nature* **551**, 256. <https://doi.org/10.1038/nature24485>.
- Lloyd-Price J et al. (2017b). Strains, functions and dynamics in the expanded Human Microbiome Project. *Nature* **550**, 61–66. <https://doi.org/10.1038/nature23889>.
- Lloyd-Price J et al. (2019). Multi-omics of the gut microbial ecosystem in inflammatory bowel diseases. *Nature* **569**, 655–662. <https://doi.org/10.1038/s41586-019-1237-9>.
- Louca S, Doebeli M (2017). Taxonomic variability and functional stability in microbial communities infected by phages. *Environmental Microbiology* **19**, 3863–3878. <https://doi.org/10.1111/1462-2920.13743>.
- Louca S, Doebeli M, Parfrey LW (2018). Correcting for 16S rRNA gene copy numbers in microbiome surveys remains an unsolved problem. *Microbiome* **6**, 41. <https://doi.org/10.1186/s40168-018-0420-9>.
- Louca S, Polz MF, et al. (2018). Function and functional redundancy in microbial systems. *Nature Ecology and Evolution* **2**, 936–943. <https://doi.org/10.1038/s41559-018-0519-1>.
- Louca S et al. (2016). Decoupling function and taxonomy in the global ocean microbiome. *Science (New York, N.Y.)* **353**, 1272–1277. <https://doi.org/10.1126/science.aaf4507>.
- Love MI et al. (2014). Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biology* **15**, 550. <https://doi.org/10.1186/s13059-014-0550-8>.
- Lozupone C, Knight R (2005). UniFrac: a New Phylogenetic Method for Comparing Microbial Communities. *Applied and environmental microbiology* **71**, 8228–8235. <https://doi.org/10.1128/AEM.71.12.8228>.
- Lozupone CA et al. (2008). The convergence of carbohydrate active gene repertoires in human gut microbes. *Proceedings of the National Academy of Sciences of the United States of America* **105**, 15076–15081. <https://doi.org/10.1073/pnas.0807339105>.
- Lu J et al. (2017). Bracken: Estimating species abundance in metagenomics data. *PeerJ Computer Science* **3**, e104. <https://doi.org/10.7717/peerj-cs.104>.
- Maghini DG et al. (2020). Improved high-molecular-weight DNA extraction, nanopore sequencing and metagenomic assembly from the human gut microbiome. *Nature Protocols* **16**, 458–471. <https://doi.org/10.1038/s41596-020-00424-x>.
- Makarova KS et al. (2015). Archaeal clusters of orthologous genes (arCOGs): An update and application for analysis of shared features between thermococcales, methanococcales, and methanobacteriales. *Life* **5**, 818–840. <https://doi.org/10.3390/life5010818>.
- Mandal S et al. (2015). Analysis of composition of microbiomes: a novel method for studying microbial composition. *Microbial Ecology in Health & Disease* **26**, 27663. <https://doi.org/10.3402/mehd.v26.27663>.
- Mandel M (1966). Deoxyribonucleic Acid Base Composition in the Genus *Pseudomonas*. *Journal of General Microbiology* **43**, 273–292. <https://doi.org/10.1099/00221287-43-2-273>.
- Manor O, Borenstein E (2017a). Revised computational metagenomic processing uncovers hidden and biologically meaningful functional variation in the human microbiome. *Microbiome* **5**, 19. <https://doi.org/10.1186/s40168-017-0231-4>.

- Manor O, Borenstein E (2017b). Systematic Characterization and Analysis of the Taxonomic Drivers of Functional Shifts in the Human Microbiome. *Cell Host and Microbe* **21**, 254–267. <https://doi.org/10.1016/j.chom.2016.12.014>.
- Martin BD et al. (2020). Modeling microbial abundances and dysbiosis with beta-binomial regression. *Annals of Applied Statistics* **14**, 94–115. <https://doi.org/10.1214/19-AOAS1283>.
- Martiny AC et al. (2013). Phylogenetic conservatism of functional traits in microorganisms. *The ISME Journal* **7**, 830–838. <https://doi.org/10.1038/ismej.2012.160>.
- Martiny AC (2019). High proportions of bacteria are culturable across major biomes. *ISME Journal* **13**, 2125–2128. <https://doi.org/10.1038/s41396-019-0410-3>.
- Maukonen J et al. (2012). The currently used commercial DNA-extraction methods give different results of clostridial and actinobacterial populations derived from human fecal samples. *FEMS Microbiology Ecology* **79**, 697–708. <https://doi.org/10.1111/j.1574-6941.2011.01257.x>.
- McCarthy A (2010). Third generation DNA sequencing: Pacific biosciences' single molecule real time technology. *Chemistry and Biology* **17**, 675–676. <https://doi.org/10.1016/j.chembiol.2010.07.004>.
- McIntyre AB et al. (2017). Comprehensive benchmarking and ensemble approaches for metagenomic classifiers. *Genome Biology* **18**, 182. <https://doi.org/10.1186/s13059-017-1299-7>.
- McLaren MR et al. (2019). Consistent and correctable bias in metagenomic sequencing experiments. *eLife* **8**, e46923. <https://doi.org/10.7554/eLife.46923>.
- McMahon K (2015). 'Metagenomics 2.0'. *Environmental Microbiology Reports* **7**, 38–39. <https://doi.org/10.1111/1758-2229.12253>.
- McNally CP et al. (2018). BURRITO: An Interactive Multi-Omic Tool for Visualizing Taxa-Function Relationships in Microbiome Data. *Frontiers in Microbiology* **9**, 365. <https://doi.org/10.3389/fmicb.2018.00365>.
- Menzel P et al. (2016). Fast and sensitive taxonomic classification for metagenomics with Kaiju. *Nature Communications* **7**, 11257. <https://doi.org/10.1038/ncomms11257>.
- Meyer F et al. (2008). The metagenomics RAST server - A public resource for the automatic phylogenetic and functional analysis of metagenomes. *BMC Bioinformatics* **9**, 386. <https://doi.org/10.1186/1471-2105-9-386>.
- Meyer F et al. (2009). FIGfams: Yet another set of protein families. *Nucleic Acids Research* **37**, 6643–6654. <https://doi.org/10.1093/nar/gkp698>.
- Mikheyev AS, Tin MMY (2014). A first look at the Oxford Nanopore MinION sequencer. *Molecular Ecology Resources* **14**, 1097–1102. <https://doi.org/10.1111/1755-0998.12324>.
- Miller CS et al. (2011). EMIRGE: Reconstruction of full-length ribosomal genes from microbial community short read sequencing data. *Genome Biology* **12**, R44. <https://doi.org/10.1186/gb-2011-12-5-r44>.
- Miossec MJ et al. (2020). Evaluation of computational methods for human microbiome analysis using simulated data. *PeerJ* **8**, e9688. <https://doi.org/10.7717/peerj.9688>.
- Morgan XC et al. (2012). Dysfunction of the intestinal microbiome in inflammatory bowel disease and treatment. *Genome Biol* **13**, R79. <https://doi.org/10.1186/gb-2012-13-9-r79>.
- Moriya Y et al. (2007). KAAS: An automatic genome annotation and pathway reconstruction server. *Nucleic Acids Research* **35**, W182–W185. <https://doi.org/10.1093/nar/gkm321>.
- Morris JJ et al. (2012). The Black Queen Hypothesis: Evolution of Dependencies through Adaptive Gene Loss. *mBio* **3**, e00036–12. <https://doi.org/10.1128/mBio.00036-12>. Copyright.
- Morton JT et al. (2017). Balance Trees Reveal Microbial Niche Differentiation. *mSystems* **2**, e00162–16. <https://doi.org/10.1128/mSystems.00162-16>.
- Morton JT et al. (2019). Establishing microbial composition measurement standards with reference frames. *Nature Communications* **10**, 2719. <https://doi.org/10.1038/s41467-019-10656-5>.
- Mosca A et al. (2016). Gut microbiota diversity and human diseases: Should we reintroduce key predators in our ecosystem? *Frontiers in Microbiology* **7**, 455. <https://doi.org/10.3389/fmicb.2016.00455>.

- Moya A, Ferrer M (2016). *Functional Redundancy-Induced Stability of Gut Microbiota Subjected to Disturbance*. *Trends in Microbiology* **24**, 402–413. <https://doi.org/10.1016/j.tim.2016.02.002>.
- Muegge BD et al. (2011). *Diet Drives Convergence in Gut Microbiome Functions Across Mammalian Phylogeny and Within Humans*. *Science* **332**, 970–974. <https://doi.org/10.1126/science.1205407>.
- Muscogiuri G et al. (2019). *Gut microbiota: a new path to treat obesity*. *International Journal of Obesity Supplements* **9**, 10–19. <https://doi.org/10.1038/s41367-019-0011-7>.
- Mysara M et al. (2017). *Reconciliation between operational taxonomic units and species boundaries*. *FEMS Microbiology Ecology* **93**, fix029. <https://doi.org/10.1093/femsec/fix029>.
- Naeem S et al. (1998). *Transcending boundaries in biodiversity research*. *Trends in Ecology and Evolution* **13**, 134–135. [https://doi.org/10.1016/s0169-5347\(97\)01316-5](https://doi.org/10.1016/s0169-5347(97)01316-5).
- Narayan NR et al. (2020). *Piphillin predicts metagenomic composition and dynamics from DADA2-corrected 16S rDNA sequences*. *BMC genomics* **21**, 56. <https://doi.org/10.1186/s12864-019-6427-1>.
- Nasko DJ et al. (2018). *RefSeq database growth influences the accuracy of k-mer-based species identification*. *Genome Biology* **19**, 156. <https://doi.org/10.1101/304972>.
- Nearing JT et al. (2018). *Denoising the Denoisers: An independent evaluation of microbiome sequence error-correction approaches*. *PeerJ* **2018**, e5364. <https://doi.org/10.7717/peerj.5364>.
- Nejman D et al. (2020). *The human tumor microbiome is composed of tumor type-specific intracellular bacteria*. *Science* **980**, 973–980. <https://doi.org/10.1126/science.aay9189>.
- NIH (2019). *A review of 10 years of human microbiome research activities at the US National Institutes of Health, Fiscal Years 2007-2016*. *Microbiome* **7**, 31. <https://doi.org/10.1186/s40168-019-0620-y>.
- Ning J, Beiko RG (2015). *Phylogenetic approaches to microbial community classification*. *Microbiome* **3**, 47. <https://doi.org/10.1186/s40168-015-0114-5>.
- Nurk S et al. (2017). *metaSPAdes: a new versatile metagenomic assembler*. *Genome Research* **27**, 824–834. <https://doi.org/10.1101/gr.213959.116>.
- Oberhardt MA et al. (2008). *Genome-scale metabolic network analysis of the opportunistic pathogen Pseudomonas aeruginosa PAO1*. *Journal of Bacteriology* **190**, 2790–2803. <https://doi.org/10.1128/JB.01583-07>.
- Ogata H et al. (1999). *KEGG: Kyoto Encyclopedia of Genes and Genomes*. *Nucleic Acids Research* **27**, 29–34. <https://doi.org/10.1093/nar/27.1.29>.
- Ogier JC et al. (2019). *rpoB, a promising marker for analyzing the diversity of bacterial communities by amplicon sequencing*. *BMC Microbiology* **19**, 171. <https://doi.org/10.1186/s12866-019-1546-z>.
- Oh J et al. (2014). *Biogeography and individuality shape function in the human skin metagenome*. *Nature* **514**, 59–64. <https://doi.org/10.1038/nature13786>.
- Olson ND et al. (2019). *Metagenomic assembly through the lens of validation: Recent advances in assessing and improving the quality of genomes assembled from metagenomes*. *Briefings in Bioinformatics* **20**, 1140–1150. <https://doi.org/10.1093/bib/bbx098>.
- Omelchenko MV et al. (2010). *Non-homologous isofunctional enzymes: A systematic analysis of alternative solutions in enzyme evolution*. *Biology Direct* **5**, 31. <https://doi.org/10.1186/1745-6150-5-31>.
- Parks DH et al. (2015). *CheckM: assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes*. *Genome research* **25**, 1043–55. <https://doi.org/10.1101/gr.186072.114>.
- Pasolli E et al. (2016). *Machine learning meta-analysis of large metagenomic datasets: tools and biological insights*. *PLoS Computational Biology* **12**, e1004977. <https://doi.org/10.1371/journal.pcbi.1004977>.
- Paulson JN et al. (2013). *Differential abundance analysis for microbial marker-gene surveys*. *Nature Methods* **10**, 1200–1202. <https://doi.org/10.1038/nmeth.2658>.

- Pawlowsky-Glahn V, Egozcue JJ (2011). Exploring compositional data with the CoDa-dendrogram. *Austrian Journal of Statistics* **40**, 103–113. <https://doi.org/10.17713/ajs.v40i1&2.202>.
- Pericard P et al. (2018). MATAM: Reconstruction of phylogenetic marker genes from short sequencing reads in metagenomes. *Bioinformatics* **34**, 585–591. <https://doi.org/10.1093/bioinformatics/btx644>.
- Pollock J et al. (2018). The Madness of Microbiome: Attempting To Find Consensus “Best Practice” for 16S Microbiome Studies. *Applied and Environmental Microbiology* **84**. Ed. by Shuang-Jiang Liu, e02627–17. <https://doi.org/10.1128/AEM.02627-17>.
- Popa O, Dagan T (2011). Trends and barriers to lateral gene transfer in prokaryotes. *Current Opinion in Microbiology* **14**, 615–623. <https://doi.org/10.1016/j.mib.2011.07.027>.
- Prakash T, Taylor TD (2012). Functional assignment of metagenomic data: challenges and applications. *Briefings in Bioinformatics* **13**, 711–727. <https://doi.org/10.1093/bib/bbs033>.
- Prodan A et al. (2020). Comparing bioinformatic pipelines for microbial 16S rRNA amplicon sequencing. *PLoS ONE* **15**, e0227434. <https://doi.org/10.1371/journal.pone.0227434>.
- Punta M et al. (2012). The Pfam protein families database. *Nucleic Acids Research* **40**, D290–D301. <https://doi.org/10.1093/nar/gkr1065>.
- Quast C et al. (2013). The SILVA ribosomal RNA gene database project: Improved data processing and web-based tools. *Nucleic Acids Research* **41**, 590–596. <https://doi.org/10.1093/nar/gks1219>.
- Rahman SF et al. (2018). Machine Learning Leveraging Genomes from Metagenomes Identifies Influential Antibiotic Resistance Genes in the Infant Gut Microbiome. *mSystems* **3**, e00123–17. <https://doi.org/10.1128/mSystems.00123-17>.
- Rasko DA et al. (2008). The pangenome structure of *Escherichia coli*: Comparative genomic analysis of *E. coli* commensal and pathogenic isolates. *Journal of Bacteriology* **190**, 6881–6893. <https://doi.org/10.1128/JB.00619-08>.
- Riley M (1993). Functions of the gene products of *Escherichia coli*. *Microbiological Reviews* **57**, 862–952. <https://doi.org/10.1128/mbr.57.4.862-952.1993>.
- Salonen A et al. (2010). Comparative analysis of fecal DNA extraction methods with phylogenetic microarray: Effective recovery of bacterial and archaeal DNA using mechanical cell lysis. *Journal of Microbiological Methods* **81**, 127–134. <https://doi.org/10.1016/j.mimet.2010.02.007>.
- Schloss PD (2020). Reintroducing *mothur*: 10 Years Later. *Applied and Environmental Microbiology* **86**, e02343–19. <https://doi.org/10.1128/aem.02343-19>.
- Schoch CL et al. (2012). Nuclear ribosomal internal transcribed spacer (ITS) region as a universal DNA barcode marker for Fungi. *Proceedings of the National Academy of Sciences of the United States of America* **109**, 6241–6246. <https://doi.org/10.1073/pnas.1117018109>.
- Schwager E et al. (2017). A Bayesian method for detecting pairwise associations in compositional data. *PLoS Computational Biology* **13**, e1005852. <https://doi.org/10.1371/journal.pcbi.1005852>.
- Segata N et al. (2011). Metagenomic biomarker discovery and explanation. *Genome Biology* **12**, R60. <https://doi.org/10.1186/gb-2011-12-6-r60>.
- Shade A (2017). Diversity is the question, not the answer. *ISME Journal* **11**, 1–6. <https://doi.org/10.1038/ismej.2016.118>.
- Shapley LS (1953). A value for n -person games. In: *Contributions to the Theory of Games*, 2. Ed. by H W Kuhn and W Tucker. Princeton, NJ: Princeton University Press, pp. 307–317. <https://doi.org/10.7249/P0295>.
- Silverman JD et al. (2017). A phylogenetic transform enhances analysis of compositional microbiota data. *eLife* **6**, e21887. <https://doi.org/10.7554/eLife.21887>.
- Sinha R et al. (2017). Assessment of variation in microbial community amplicon sequencing by the Microbiome Quality Control (MBQC) project consortium. *Nature Biotechnology* **35**, 1077–1086. <https://doi.org/10.1038/nbt.3981>.
- Spang A et al. (2015). Complex archaea that bridge the gap between prokaryotes and eukaryotes. *Nature* **521**, 173–179. <https://doi.org/10.1038/nature14447>.

- Spencer SJ et al. (2016). Massively parallel sequencing of single cells by epicPCR links functional genes with phylogenetic markers. *ISME Journal* **10**, 427–436. <https://doi.org/10.1038/ismej.2015.124>.
- Sperling JL et al. (2017). Comparison of bacterial 16S rRNA variable regions for microbiome surveys of ticks. *Ticks and Tick-borne Diseases* **8**, 453–461. <https://doi.org/10.1016/j.ttbdis.2017.02.002>.
- 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. *International Journal of Systematic Bacteriology* **44**, 846–849. <https://doi.org/10.1099/00207713-44-4-846>.
- Staley J, Konopka A (1985). Measurement of In Situ Activities of Nonphotosynthetic Microorganisms in Aquatic and Terrestrial Habitats. *Annual Review of Microbiology* **39**, 321–346. <https://doi.org/10.1146/annurev.micro.39.1.321>.
- Stein JL et al. (1996). Characterization of uncultivated prokaryotes: Isolation and analysis of a 40-kilobase-pair genome fragment from a planktonic marine archaeon. *Journal of Bacteriology* **178**, 591–599. <https://doi.org/10.1128/jb.178.3.591-599.1996>.
- Steinberger M, Söding J (2017). MMseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nature Biotechnology* **35**, 1026–1028. <https://doi.org/10.1038/nbt.3988>.
- Steinberger M, Söding J (2018). Clustering huge protein sequence sets in linear time. *Nature Communications* **9**, 2542. <https://doi.org/10.1038/s41467-018-04964-5>.
- Suardana IW (2014). Analysis of Nucleotide Sequences of the 16S rRNA Gene of Novel *Escherichia coli* Strains Isolated from Feces of Human and Bali Cattle. *Journal of Nucleic Acids* **2014**, 475754. <https://doi.org/10.1155/2014/475754>.
- Sun DL et al. (2013). Intragenomic heterogeneity of 16S rRNA genes causes overestimation of prokaryotic diversity. *Applied and Environmental Microbiology* **79**, 5962–5969. <https://doi.org/10.1128/AEM.01282-13>.
- Sun S et al. (2020). Inference-based accuracy of metagenome prediction tools varies across sample types and functional categories. *Microbiome* **8**, 46. <https://doi.org/10.1186/s40168-020-00815-y>.
- Sunagawa S et al. (2015). Structure and function of the global ocean microbiome. *Science* **348**, 1261359. <https://doi.org/10.1126/science.1261359>.
- Sze MA, Schloss PD (2016). Looking for a signal in the noise: Revisiting obesity and the microbiome. *mBio* **7**, e01018–16. <https://doi.org/10.1128/mBio.01018-16>.
- Tatusov RL et al. (2000). The COG database: a tool for genome-scale analysis of protein functions and evolution. *Nucleic Acids Research* **28**, 33–36. <https://doi.org/10.1093/nar/28.1.33>.
- Tessler M et al. (2017). Large-scale differences in microbial biodiversity discovery between 16S amplicon and shotgun sequencing. *Scientific Reports* **7**, 6589. <https://doi.org/10.1038/s41598-017-06665-3>.
- Tettelin H et al. (2005). Genome analysis of multiple pathogenic isolates of *Streptococcus agalactiae*: Implications for the microbial 'pan-genome'. *Proceedings of the National Academy of Sciences of the United States of America* **102**, 3950–3955. <https://doi.org/10.1073/pnas.0508532102>.
- Thomas PD et al. (2007). Ontology annotation: mapping genomic regions to biological function. *Current Opinion in Chemical Biology* **11**, 4–11. <https://doi.org/10.1016/j.cbpa.2006.11.039>.
- Thompson LR et al. (2017). A communal catalogue reveals Earth's multiscale microbial diversity. *Nature* **551**, 457–463. <https://doi.org/10.1038/nature24621>.
- Thorsen J et al. (2016). Large-scale benchmarking reveals false discoveries and count transformation sensitivity in 16S rRNA gene amplicon data analysis methods used in microbiome studies. *Microbiome* **4**, 62. <https://doi.org/10.1186/s40168-016-0208-8>.
- Treem W et al. (1994). Fecal Short-Chain Fatty Acids in Children with Inflammatory Bowel Disease. *Journal of Pediatric Gastroenterology and Nutrition* **18**, 159–164. <https://doi.org/10.1097/00005176-199402000-00007>.

- Truong DT et al. (2015). *MetaPhlan2 for enhanced metagenomic taxonomic profiling*. *Nature Methods* **12**, 902–903. <https://doi.org/10.1038/nmeth.3589>.
- Turnbaugh PJ et al. (2008). *Diet-Induced Obesity Is Linked to Marked but Reversible Alterations in the Mouse Distal Gut Microbiome*. *Cell Host and Microbe* **3**, 213–223. <https://doi.org/10.1016/j.chom.2008.02.015>.
- Turnbaugh PJ et al. (2009). *A core gut microbiome in obese and lean twins*. *Nature* **457**, 480–484. <https://doi.org/10.1038/nature07540>.
- Vandeputte D et al. (2017). *Quantitative microbiome profiling links gut community variation to microbial load*. *Nature* **551**, 507–511. <https://doi.org/10.1038/nature24460>.
- Vatanen T et al. (2018). *Genomic variation and strain-specific functional adaptation in the human gut microbiome during early life*. *Nature Microbiology* **4**, 470–479. <https://doi.org/10.1038/s41564-018-0321-5>.
- Venegas DP et al. (2019). *Short chain fatty acids (SCFAs) mediated gut epithelial and immune regulation and its relevance for inflammatory bowel diseases*. *Frontiers in Immunology* **10**, 277. <https://doi.org/10.3389/fimmu.2019.00277>.
- Venter JC et al. (2004). *Environmental Genome Shotgun Sequencing of the Sargasso Sea*. *Science* **304**, 66–74. <https://doi.org/10.1126/science.1093857>.
- Verster AJ, Borenstein E (2018). *Competitive lottery-based assembly of selected clades in the human gut microbiome*. *Microbiome* **6**, 186. <https://doi.org/10.1186/s40168-018-0571-8>.
- Větrovský T, Baldrian P (2013). *The Variability of the 16S rRNA Gene in Bacterial Genomes and Its Consequences for Bacterial Community Analyses*. *PLoS ONE* **8**, e57923. <https://doi.org/10.1371/journal.pone.0057923>.
- Vincent C et al. (2013). *Reductions in intestinal Clostridiales precede the development of nosocomial Clostridium difficile infection*. *Microbiome* **1**, 18. <https://doi.org/10.1186/2049-2618-1-18>.
- Wang Y, Qian PY (2014). *Conserved Regions in 16S Ribosome RNA Sequences and Primer Design for Studies of Environmental Microbes*. *Encyclopedia of Metagenomics*. https://doi.org/10.1007/978-1-4614-6418-1_772-1.
- Watson EJ et al. (2019). *Human faecal collection methods demonstrate a bias in microbiome composition by cell wall structure*. *Scientific Reports* **9**, 16831. <https://doi.org/10.1038/s41598-019-53183-5>.
- Watterson WJ et al. (2020). *Droplet-based high-throughput cultivation for accurate screening of antibiotic resistant gut microbes*. *eLife* **9**, e56998. <https://doi.org/10.7554/eLife.56998>.
- Weiss S et al. (2017). *Normalization and microbial differential abundance strategies depend upon data characteristics*. *Microbiome* **5**, 27. <https://doi.org/10.1186/s40168-017-0237-y>.
- Welch RA et al. (2002). *Extensive mosaic structure revealed by the complete genome sequence of uropathogenic Escherichia coli*. *Proceedings of the National Academy of Sciences of the United States of America* **99**, 17020–4. <https://doi.org/10.1073/pnas.252529799>.
- Wemheuer F et al. (2020). *Tax4Fun2: prediction of habitat-specific functional profiles and functional redundancy based on 16S rRNA gene sequences*. *Environmental Microbiome* **15**, 11. <https://doi.org/10.1101/490037>.
- Wheeler NE et al. (2016). *A profile-based method for identifying functional divergence of orthologous genes in bacterial genomes*. *Bioinformatics* **32**, 3566–3574. <https://doi.org/10.1093/bioinformatics/btw518>.
- Willis AD (2019). *Rigorous Statistical Methods for Rigorous Microbiome Science*. *mSystems* **4**, e00117–19. <https://doi.org/10.1128/msystems.00117-19>.
- Willis C et al. (2019). *Influence of 16S rRNA variable region on perceived diversity of marine microbial communities of the Northern North Atlantic*. *FEMS Microbiology Letters* **366**, fnz152. <https://doi.org/10.1093/femsle/fnz152>.
- Wilson GA et al. (2005). *Orphans as taxonomically restricted and ecologically important genes*. *Microbiology* **151**, 2499–2501. <https://doi.org/10.1099/mic.0.28146-0>.
- Wirbel J et al. (2019). *Meta-analysis of fecal metagenomes reveals global microbial signatures that are specific for colorectal cancer*. *Nature Medicine* **25**, 679–689. <https://doi.org/10.1038/s41591-019-0406-6>.

- Woese CR (1987). *Bacterial evolution*. *Microbiological Reviews* **51**, 221–271. <https://doi.org/10.1128/mbr.51.2.221-271.1987>.
- Woese CR et al. (1980). *Secondary structure model for bacterial 16S ribosomal RNA: Phylogenetic, enzymatic and chemical evidence*. *Nucleic Acids Research* **8**, 2275–2294. <https://doi.org/10.1093/nar/8.10.2275>.
- Woese CR, Fox GE (1977). *Phylogenetic structure of the prokaryotic domain: The primary kingdoms*. *Proceedings of the National Academy of Sciences of the United States of America* **74**, 5088–5090. <https://doi.org/10.1073/pnas.74.11.5088>.
- Wood DE et al. (2019). *Improved metagenomic analysis with Kraken 2*. *Genome Biology* **20**, 257. <https://doi.org/10.1186/s13059-019-1891-0>.
- Wright EK et al. (2015). *Recent advances in characterizing the gastrointestinal microbiome in Crohn's disease: a systematic review*. *Inflamm Bowel Dis* **21**, 1219–1228. <https://doi.org/10.1097/MIB.0000000000000382>.
- Wu D et al. (2013). *Systematic Identification of Gene Families for Use as "Markers" for Phylogenetic and Phylogeny-Driven Ecological Studies of Bacteria and Archaea and Their Major Subgroups*. *PLoS ONE* **8**, e77033. <https://doi.org/10.1371/journal.pone.0077033>.
- Xu Y, Zhao F (2018). *Single-cell metagenomics: challenges and applications*. *Protein and Cell* **9**, 501–510. <https://doi.org/10.1007/s13238-018-0544-5>.
- Ye SH et al. (2019). *Benchmarking Metagenomics Tools for Taxonomic Classification*. *Cell* **178**, 779–794. <https://doi.org/10.1016/j.cell.2019.07.010>.
- Ye Y, Doak TG (2011). *A Parsimony Approach to Biological Pathway Reconstruction/Inference for Metagenomes*. *PLoS Computational Biology* **5**, e1000465. <https://doi.org/10.1002/9781118010518.ch52>.
- Zaneveld JR et al. (2010). *Ribosomal RNA diversity predicts genome diversity in gut bacteria and their relatives*. *Nucleic Acids Research* **38**, 3869–3879. <https://doi.org/10.1093/nar/gkq066>.
- Zaneveld JR, Thurber RL (2014). *Hidden state prediction: A modification of classic ancestral state reconstruction algorithms helps unravel complex symbioses*. *Frontiers in Microbiology* **5**, 431. <https://doi.org/10.3389/fmicb.2014.00431>.
- Zemb O et al. (2020). *Absolute quantitation of microbes using 16S rRNA gene metabarcoding: A rapid normalization of relative abundances by quantitative PCR targeting a 16S rRNA gene spike-in standard*. *MicrobiologyOpen* **9**, e977. <https://doi.org/10.1002/mbo3.977>.
- Zhou J et al. (2015). *High-Throughput Metagenomic Technologies for Complex Microbial Community Analysis: Open and Closed Formats*. *mBio* **6**, e02288–14. <https://doi.org/10.1128/mBio.02288-14>. Copyright.
- Zhou W et al. (2019). *Longitudinal multi-omics of host-microbe dynamics in prediabetes*. *Nature* **569**, 663–671. <https://doi.org/10.1038/s41586-019-1236-x>.
- Zhou YH, Gallins P (2019). *A review and tutorial of machine learning methods for microbiome host trait prediction*. *Frontiers in Genetics* **10**, 579. <https://doi.org/10.3389/fgene.2019.00579>.
- Zuckermandl E, Pauling L (1965). *Molecules as documents of history*. *Journal of Theoretical Biology* **8**, 357–366. [https://doi.org/10.1016/0022-5193\(65\)90083-4](https://doi.org/10.1016/0022-5193(65)90083-4).