

Abstract

Alterations in the composition of the microbiota have been implicated in many diseases. The Human Microbiome Project (HMP) provides a comprehensive reference dataset of the “normal” human microbiome of 242 healthy adults at five major body sites. The HMP used both 16S ribosomal RNA gene sequencing and whole-genome metagenomic sequencing to profile the subjects’ microbial communities. However, accessing and analyzing the HMP dataset still presents technical and bioinformatic challenges, as researchers must import the microbiome data, integrate phylogenetic trees, and access and merge public and restricted metadata. In this issue, the *HMP16SData* R/Bioconductor package developed by Schiffer and colleagues (Am J Epidemiol. XXX; XX (XX): XX–XXX) greatly simplifies access to the HMP data by combining 16S taxonomic abundance data, public patient metadata, and phylogenetic trees as a single data object. The authors also provide an interface for users with approved dbGaP projects to easily retrieve and merge the controlled-access HMP metadata. This package has a broad range of appeal to researchers across disciplines and with various levels of expertise in using R and/or other statistical tools. This will translate to improved data accessibility for public health research, with data from healthy individuals serving as a reference for disease-associated studies.

Keywords: Bioconductor, databases, Human Microbiome Project, microbiome, R

1
2
3 *Large datasets such as the Human Microbiome Project are needed for robust epidemiology of*
4 *microbiome-associated diseases*
5
6
7

8
9 The human microbiota comprises the bacterial, fungal, archaeal and viral occupants of the
10 human body. An increasing number of health conditions have been linked to the
11 gastrointestinal (e.g. inflammatory bowel disease (1), colorectal cancer (2), obesity (3), type 2
12 diabetes (3), and rheumatoid arthritis (4)), oral (e.g. periodontitis (5), pancreatic cancer (6)),
13 skin (e.g. dermatitis (7) and other cutaneous diseases (8)), and vaginal (e.g. bacterial vaginosis
14 (9)) microbiota. In order to understand how the microbiome changes in disease, it is crucial to
15 first understand the microbial composition and variance within healthy populations. MetaHIT
16 (10) was the first large-scale survey of healthy adult stool microbiota, while the first phase of
17 the Human Microbiome Project (HMP) (11) provided a comprehensive reference dataset of the
18 “normal” human microbiome of healthy individuals at five major body sites (GI tract, nasal,
19 skin, oral cavity, and vagina). Applying these data to create well-designed epidemiological
20 studies and meta-analyses is an important step in determining how dysbiosis (deviation from
21 the normal microbiome) in microbial composition and functional profiles may contribute to
22 disease etiology. While the HMP data is freely available to researchers at the HMP Data
23 Analysis and Coordination Center (HMPDACC) website (12), its format is not readily
24 importable into widely-used statistical software such as R or Stata, presenting a significant
25 technical challenge to researchers. Schiffer and colleagues have substantially reduced these
26 technical challenges by creating the *HMP16SData* R/Bioconductor package (13) which greatly
27 simplifies researcher access to HMP dataset. The *HMP16SData* removes technical hurdles to
28 data processing and provides access to the HMP 16S rRNA data as an object that is ready to
29 analyze. This will provide epidemiologists with a powerful tool to help to identify robust
30 epidemiological links between the microbiome and disease.
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 *The HMP data from HMPDACC is a valuable but technically challenging epidemiological*
4 *resource*
5
6

7
8
9 Between 2008 and 2010, the HMP phase I collected thousands of samples from 15 (male) or
10
11 18 (female) distinct body sites of 242 healthy adults between the ages of 18 and 40 years over
12
13 multiple time points. To characterize the composition of samples' microbial communities, the
14
15 454 FLX Titanium platform was used to sequence the V35 hypervariable region of the 16S
16
17 ribosomal RNA (rRNA) gene. The V13 hypervariable region of the 16S rRNA gene was also
18
19 sequenced for a subset of the samples for a more complete picture of a community (14) and to
20
21 allow for methods comparison. Finally, a subset of samples underwent whole-genome shotgun
22
23 (WGS) or metagenomic sequencing of the whole community DNA using the Illumina GAIIx
24
25 platform. Although 16S profiling is currently the most commonly used and most cost-effective
26
27 tool for microbiome analysis, WGS (which uses sequencing with random primers to sequence
28
29 overlapping regions of a genome) allows for more accurate definition of taxa at the species
30
31 level, as well as allowing for identification and quantitation of functional genes.
32
33
34
35
36

37
38 A wealth of HMP data is available for download at the HMP's Data Analysis Coordination
39
40 Center (HMPDACC) (12): 16S rRNA data for hypervariable regions V13 and V35 (e.g. OTU
41
42 tables describing the abundance of taxa in samples, phylogenetic trees describing the
43
44 relationships of OTUs, taxonomy data, WGS sequencing data, and some accompanying sample
45
46 metadata such as sample type, patient sex, sequencing center, and visit number. Sensitive
47
48 metadata such as subject age and medical history are kept confidential, but researchers who
49
50 submit an application can be granted access through the National Center for Biotechnology
51
52 Information (NCBI) database of Genotypes and Phenotypes (dbGAP). Despite its aim of
53
54 having accessible data and providing user-friendly data retrieval, retrieval and downstream
55
56 analysis of HMP data from the DACC in its current format still presents as a considerable
57
58
59
60

1
2
3 bioinformatic challenge, especially to researchers with limited knowledge of HMPDACC and
4 dbGaP procedures. Researchers are required to import the microbiome taxonomy and
5 abundance data, integrate phylogenetic trees, and correctly merge metadata from hundreds of
6 patients with thousands of samples. Furthermore, using restricted metadata requires an
7 application to dbGaP, after which the clinical metadata must be downloaded and merged.
8 Substantial expertise is needed to successfully merge these large datasets into suitable formats
9 for downstream analysis. In order to address this issue, the authors of the *phyloseq*
10 R/Bioconductor package (15) have previously provided a guide on how to import the HMPv35
11 dataset from HMPDACC (16) and have made the 16S rRNA data from V35 hypervariable
12 regions publicly available as a processed single R object (17). This is a very useful resource,
13 but it does not include the V13 or WGS datasets, and it contains only the limited, publicly-
14 available patient metadata. Here, Schiffer *et al* have improved on this by creating
15 *HMP16SData*, an R dataset that integrates the HMPv13 and HMPv35 16S rRNA data and
16 provides easy access of restricted patient data from dbGaP to researchers with approved dbGaP
17 projects. Furthermore, *HMP16SData* is easily merged with the HMP WGS data, which is
18 available via the *CuratedMetagenomeData* R/Bioconductor package(18).

19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42 *Ready integration of metadata, taxonomy, abundance, and phylogenetic relationships*
43 *facilitates quantitative analysis*

44
45
46
47 In 16S rRNA analysis, raw sequencing data is processed by open source computational tools
48 and pipelines such as mothur (19) and QIIME (20), which group sequences into operational
49 taxonomic units (OTUs) based on sequence similarity (usually 97% identity, corresponding
50 approximately but not exactly to species). The number of times each OTU sequence is detected
51 in each sample is summarized as an OTU table. OTU counts can be summarized at various
52 taxonomic levels (e.g., class, order, family, genus) for analysis.
53
54
55
56
57
58
59
60

1
2
3 The Bioconductor package *phyloseq* (15) in R (21) is the most popular tool for downstream
4 analysis of microbial sequencing data. It facilitates statistical analysis and creation of
5 publication-quality graphics, and enables reproducible research when used with documentation
6 tools such as markdown (22). A distinctive feature of *phyloseq* is the integration of OTU-
7 clustered data, taxonomic assignments, and associated sample data as a *phyloseq* object. This
8 allows users to investigate the relative abundance and diversity of organisms at various
9 taxonomic levels, which is especially useful in instances where analyses at taxonomic ranks
10 higher than species provide more ecologically meaningful information. Another useful feature
11 of *phyloseq* is the ability to easily agglomerate taxa by their taxonomic ranks. The data
12 integration provided by *phyloseq* greatly simplifies analysis of microbial within the R
13 environment.
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28

29 A special feature of a microbiome data matrix is that the distribution of individual OTUs is
30 highly skewed and often sparse (23). For example, the bacterial abundance in the human gut
31 microbiome consists of a high proportion of zero counts at lower taxonomic levels (24). The
32 presence of excess zeros presents a challenge when analyzing microbiome data, particularly
33 when comparing between groups. Two common approaches to solve this problem include
34 normalizing transformation (e.g., variance-stabilizing transformations and linear modelling)
35 and log transformation using a generalized linear model (GLM). Several specialized R
36 /Bioconductor packages greatly simplify this statistical modeling processes and are well-
37 integrated with the *phyloseq* package; these include *MetagenomeSeq* (23), which was
38 developed specifically for marker gene analysis, as well as RNA-Seq focused R packages such
39 as *DESeq2* (25) and *edgeR* (26).
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

55 In this issue, Schiffer *et al.* (13) have provided an elegant way for researchers to quickly access
56 and analyze the HMP phase 1 data. The authors developed the *HMP16SData* R/Bioconductor
57
58
59
60

1
2
3 package by combining HMP 16S taxonomic abundance data, public and (optionally) restricted
4 patient metadata, and phylogenetic trees as a single data object, which can then be easily
5 converted into a *phyloseq* object or alternatively exported in .csv, STATA, SAS, or SPSS
6 formats for use with other statistical software. The *HMP16SData* package is easy to install via
7 Bioconductor (27), and clear, helpful online documentation is also available in the package
8 vignette (28) as well as the authors' github site (29). In addition, after researchers have
9 requested and obtained access to dbGaP, the *HMP16SData* package includes the
10 *attach_dbGaP* function, which allows decryption and attachment of restricted patient data
11 from dbGaP into the other data. By removing barriers to data access and management, the
12 authors have made microbiome data from the HMP much more accessible to the research
13 community.

14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 *Limitations of the HMP16SData*

31
32
33 Researchers who wish to compare their own datasets to the HMP datasets should be aware of
34 inherent biases due to different sampling handling, DNA extraction protocols, 16S rRNA gene
35 primer selection, sequencing platforms and bioinformatics processing pipelines (30).
36 Sequencing of the 16S rRNA gene in the HMP samples relied on the now discontinued Roche
37 454 platform, whereas the majority of amplicon sequencing is now done on the Illumina
38 platform due to its higher throughput and lower cost. The 454 platform yields longer but fewer
39 sequences, while the Illumina provides shorter reads but at much greater sequencing depth.
40 Furthermore, the platforms are prone to different types of sequencing error, which require
41 different forms of error correction during bioinformatic processing of the sequencing data (31,
42 32). The difference in total number of sequences per sample must be considered when
43 comparing alpha diversity (total number of taxa) in populations between Illumina and 454 data.
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 However, even if data is rarefied (subsampling to the same number of reads per sample), these
4 platform differences will introduce some biases the researcher should be aware of.
5
6

7
8 In addition, the taxonomy reference databases and analysis pipelines which were used to
9 process the HMP data have had many version updates (e.g., QIIME 1.3.0 vs. 2.2018.8 current,
10 mothur 1.1.8 vs. 1.40 current); these software differences could introduce additional bias.
11 Finally, it is becoming increasingly common to process Illumina amplicon sequencing data by
12 inferring of amplicon sequence variants (e.g. exact sequence variants) rather than OTUs with
13 97% sequence similarity. Two examples of sequence variant-calling software include DADA2
14 (33) and Deblur (34); both of these are implemented in QIIME 2 (35). If the same sequencing
15 data are processed by both QIIME 1.9 and QIIME 2.0, general patterns of taxonomic
16 composition and beta diversity will be similar, but due to differences in correcting for
17 sequencing error, data produced by QIIME 2.0 will give a more conservative and more accurate
18 estimation of the alpha diversity of a sample (36).
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

35 Because there is still no batch normalization method for microbiome data, researchers carrying
36 out meta-analyses using microbiome data such as the HMP data must be aware of these batch
37 and data effects. For more comparable results across studies, certain protocol choices including
38 using relative (not absolute) diversity measures, phylogenetic (not taxonomic) analyses, and
39 quantitative (not presence or absence) measures may buffer the inter-protocol variations (30,
40 36). Schiffer *et al* recognize the biases that occur with cross-study comparisons, and
41 acknowledge the need to reprocess raw HMP data with modern bioinformatics tools.
42
43
44
45
46
47
48
49
50

51 *Conclusions*

52
53
54 In summary, the *HMP16SData* R/Bioconductor package developed by Schiffer and colleagues
55 (13) greatly simplifies access to the Human Microbiome Project data, a landmark data
56 resource. In our opinion, this package has a broad range of appeal to researchers across
57
58
59
60

1
2
3 disciplines and with various levels of expertise in using R and/or other statistical tools. This
4
5 will translate to improved data accessibility for public health research, with data from healthy
6
7 individuals serving as a reference for disease-associated studies, or as a baseline for comparing
8
9 the Western population with microbiomes from other geographic and ethnic cohorts. As
10
11 additional large-scale datasets become available, *HMP16SData* will be an invaluable tool to
12
13 the research community to easily access and analyze these resulting datasets and better
14
15 understand host-microbe interactions.
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

References

1. Morgan XC, Tickle TL, Sokol H, et al. Dysfunction of the intestinal microbiome in inflammatory bowel disease and treatment. *Genome Biol* 2012;13(9):R79.
2. Louis P, Hold GL, Flint HJ. The gut microbiota, bacterial metabolites and colorectal cancer. *Nat Rev Microbiol* 2014;12(10):661-72.
3. Hartstra AV, Bouter KE, Backhed F, et al. Insights into the role of the microbiome in obesity and type 2 diabetes. *Diabetes Care* 2015;38(1):159-65.
4. Scher JU, Abramson SB. The microbiome and rheumatoid arthritis. *Nat Rev Rheumatol* 2011;7(10):569-78.
5. Jorth P, Turner KH, Gumus P, et al. Metatranscriptomics of the human oral microbiome during health and disease. *MBio* 2014;5(2):e01012-14.
6. Fan X, Alekseyenko AV, Wu J, et al. Human oral microbiome and prospective risk for pancreatic cancer: a population-based nested case-control study. *Gut* 2018;67(1):120-7.
7. Kong HH, Oh J, Deming C, et al. Temporal shifts in the skin microbiome associated with disease flares and treatment in children with atopic dermatitis. *Genome Res* 2012;22(5):850-9.
8. Findley K, Grice EA. The skin microbiome: a focus on pathogens and their association with skin disease. *PLoS Pathog* 2014;10(10):e1004436.
9. Ma B, Forney LJ, Ravel J. Vaginal microbiome: rethinking health and disease. *Annu Rev Microbiol* 2012;66:371-89.
10. Qin J, Li R, Raes J, et al. A human gut microbial gene catalogue established by metagenomic sequencing. *Nature* 2010;464(7285):59-65.
11. The Human Microbiome Project Consortium. Structure, function and diversity of the healthy human microbiome. *Nature* 2012;486(7402):207-14.
12. NIH Human Microbiome Project. The HMP Data Analysis and Coordination Center. (<https://hmpdacc.org/hmp/>). (Accessed 2 Nov 2018).
13. Schiffer L, Azhar R, Shepherd L, et al. HMP16SData: Efficient Access to the Human Microbiome Project through Bioconductor. *American Journal of Epidemiology* 2018.
14. Human Microbiome Project Jumpstart Consortium. Evaluation of 16S rDNA-based community profiling for human microbiome research. *PLoS ONE* 2012;7(6):e39315.
15. McMurdie PJ, Holmes S. phyloseq: an R package for reproducible interactive analysis and graphics of microbiome census data. *PLoS One* 2013;8(4):e61217.
16. McMurdie PJ, Holmes S. Human Microbiome Project. (http://joey711.github.io/phyloseq-demo/HMP_import_example.html). (Accessed 2 Nov 2018).
17. McMurdie PJ, Holmes S. Already-imported HMPv35 dataset. (<http://joey711.github.io/phyloseq-demo/HMPv35.RData>). (Accessed 2 Nov 2018).
18. Pasolli E, Schiffer L, Manghi P, et al. Accessible, curated metagenomic data through ExperimentHub. *Nat Methods* 2017;14(11):1023-4.
19. Schloss PD, Westcott SL, Ryabin T, et al. Introducing mothur: open-source, platform-independent, community-supported software for describing and comparing microbial communities. *Applied and environmental microbiology* 2009;75(23):7537-41.
20. Caporaso JG, Kuczynski J, Stombaugh J, et al. QIIME allows analysis of high-throughput community sequencing data. *Nat Methods* 2010;7(5):335-6.

- 1
 - 2
 - 3
 - 4
 - 5
 - 6
 - 7
 - 8
 - 9
 - 10
 - 11
 - 12
 - 13
 - 14
 - 15
 - 16
 - 17
 - 18
 - 19
 - 20
 - 21
 - 22
 - 23
 - 24
 - 25
 - 26
 - 27
 - 28
 - 29
 - 30
 - 31
 - 32
 - 33
 - 34
 - 35
 - 36
 - 37
 - 38
 - 39
 - 40
 - 41
 - 42
 - 43
 - 44
 - 45
 - 46
 - 47
 - 48
 - 49
 - 50
 - 51
 - 52
 - 53
 - 54
 - 55
 - 56
 - 57
 - 58
 - 59
 - 60
21. R Core Team. R: A Language and Environment for Statistical Computing. Vienna, Austria; 2018. (<https://www.R-project.org/>). (Accessed).
22. Allaire J.J., Horner J., Marti V., et al. markdown: 'Markdown' Rendering for R. 2017. (<https://CRAN.R-project.org/package=markdown>). (Accessed).
23. Paulson JN, Stine OC, Bravo HC, et al. Differential abundance analysis for microbial marker-gene surveys. *Nat Methods* 2013.
24. Xia Y, Sun J. Hypothesis Testing and Statistical Analysis of Microbiome. *Genes Dis* 2017;4(3):138-48.
25. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol* 2014;15(12):550.
26. Robinson MD, McCarthy DJ, Smyth GK. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* 2010;26(1):139-40.
27. Huber W, Carey VJ, Gentleman R, et al. Orchestrating high-throughput genomic analysis with Bioconductor. *Nat Methods* 2015;12(2):115-21.
28. Schiffer L, Azhar R, Shepherd L, et al. HMP16SData 2018. (<https://bioconductor.org/packages/devel/data/experiment/vignettes/HMP16SData/inst/doc/HMP16SData.html>). (Accessed 1 Nov 2018).
29. Schiffer L, Azhar R, Shepherd L, et al. 16S rRNA Sequencing Data from the Human Microbiome Project. (<https://github.com/waldronlab/HMP16SData>). (Accessed 2 Nov 2018).
30. Sinha R, Abu-Ali G, Vogtmann E, et al. Assessment of variation in microbial community amplicon sequencing by the Microbiome Quality Control (MBQC) project consortium. *Nat Biotechnol* 2017;35(11):1077-86.
31. Quince C, Lanzen A, Davenport RJ, et al. Removing noise from pyrosequenced amplicons. *BMC Bioinformatics* 2011;12:38.
32. Luo C, Tsementzi D, Kyrpides N, et al. Direct comparisons of Illumina vs. Roche 454 sequencing technologies on the same microbial community DNA sample. *PLoS One* 2012;7(2):e30087.
33. Callahan BJ, McMurdie PJ, Rosen MJ, et al. DADA2: High-resolution sample inference from Illumina amplicon data. *Nat Methods* 2016;13(7):581-3.
34. Amir A, McDonald D, Navas-Molina JA, et al. Deblur Rapidly Resolves Single-Nucleotide Community Sequence Patterns. *mSystems* 2017;2(2).
35. Bolyen E, Rideout JR, Dillon MR, et al. QIIME 2: Reproducible, interactive, scalable, and extensible microbiome data science. *PeerJ Preprints* 2018;6:e27295v1.
36. Nearing JT, Douglas GM, Comeau AM, et al. Denoising the Denoisers: an independent evaluation of microbiome sequence error-correction approaches. *PeerJ* 2018;6:e5364.