

University of Nebraska - Lincoln

DigitalCommons@University of Nebraska - Lincoln

---

Faculty Publications in Food Science and  
Technology

Food Science and Technology Department

---

6-14-2012

## A framework for human microbiome research

Barbara A. Methe

Jacques Izard

Follow this and additional works at: <https://digitalcommons.unl.edu/foodsciefacpub>



Part of the [Food Science Commons](#)

---

This Article is brought to you for free and open access by the Food Science and Technology Department at DigitalCommons@University of Nebraska - Lincoln. It has been accepted for inclusion in Faculty Publications in Food Science and Technology by an authorized administrator of DigitalCommons@University of Nebraska - Lincoln.

# A framework for human microbiome research

The Human Microbiome Project Consortium\*

**A variety of microbial communities and their genes (the microbiome) exist throughout the human body, with fundamental roles in human health and disease. The National Institutes of Health (NIH)-funded Human Microbiome Project Consortium has established a population-scale framework to develop metagenomic protocols, resulting in a broad range of quality-controlled resources and data including standardized methods for creating, processing and interpreting distinct types of high-throughput metagenomic data available to the scientific community. Here we present resources from a population of 242 healthy adults sampled at 15 or 18 body sites up to three times, which have generated 5,177 microbial taxonomic profiles from 16S ribosomal RNA genes and over 3.5 terabases of metagenomic sequence so far. In parallel, approximately 800 reference strains isolated from the human body have been sequenced. Collectively, these data represent the largest resource describing the abundance and variety of the human microbiome, while providing a framework for current and future studies.**

Advances in sequencing technologies coupled with new bioinformatic developments have allowed the scientific community to begin to investigate the microbes that inhabit our oceans, soils, the human body and elsewhere<sup>1</sup>. Microbes associated with the human body include eukaryotes, archaea, bacteria and viruses, with bacteria alone estimated to outnumber human cells within an individual by an order of magnitude. Our knowledge of these communities and their gene content, referred to collectively as the human microbiome, has until now been limited by a lack of population-scale data detailing their composition and function.

The US NIH-funded Human Microbiome Project Consortium (HMP) brought together a broad collection of scientific experts to explore these microbial communities and their relationships with their human hosts. As such, the HMP<sup>2</sup> has focused on producing reference genomes (viral, bacterial and eukaryotic), which provide a critical framework for subsequent metagenomic annotation and analysis, and on generating a baseline of microbial community structure and function from an adult cohort defined by a carefully delineated set of clinical inclusion and exclusion criteria that we term 'healthy' in this study (<http://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/GetPdf.cgi?id=phd002854.2>). Investigations of the microbiome from this cohort incorporated several complementary analyses including: 16S ribosomal RNA (rRNA) gene sequence (16S) and taxonomic profiles, whole-genome shotgun (WGS) or metagenomic sequencing of whole community DNA, and alignment of the assembled sequences to the reference microbial genomes from the human body<sup>3,4</sup>. Thus, the HMP complements other large-scale sequence-based human microbiome projects such as the MetaHIT project<sup>5</sup>, which focused on examination of the gut microbiome using WGS data including samples from cohorts exhibiting a wide range of health statuses and physiological characteristics.

Additional projects supported by the HMP are investigating the association of specific components and dynamics of the microbiome with a variety of disease conditions, developing tools and technology including isolating and sequencing uncultured organisms, and studying the ethical, legal and social implications of human microbiome research (<http://commonfund.nih.gov/hmp/fundedresearch.aspx>). A comprehensive list of current publications from HMP projects is available at <http://commonfund.nih.gov/hmp/publications.aspx>.

Here we detail the resources created so far by the HMP initiative including: clinical specimens (samples), reference genomes, sequencing and annotation protocols, methods and analyses. We describe the thousands of samples obtained from 15 or 18 distinct body sites from 242 donors over multiple time points that were processed at two clinical centres (Baylor College of Medicine (BCM) and Washington University School of Medicine). We also describe the laboratory and computational protocols developed for reliably generating and interpreting the human microbiome data. HMP resources include both protocols for, and the subsequent data generated from, 16S and metagenomic sequencing of human microbiome samples. During this study, these protocols were rigorously standardized and quality controlled for simultaneous use across four sequencing centres (BCM Human Genome Sequencing Center, The Broad Institute of Massachusetts Institute of Technology (MIT) and Harvard, the J. Craig Venter Institute and The Genome Institute at Washington University School of Medicine). In particular, we focus on the production of the first phase of metagenomic data sets (phase I) used for subsequent in-depth analyses, and we summarize standards and recommendations based on our experiences generating and analysing these data. An additional set of publications (many included in the references and in those of ref. 4) describe in further detail the microbial ecology and microbiological implications of these data. Collectively these resources and analyses represent an important framework for human microbiome research.

## HMP resource organization

Supplementary Fig. 1 summarizes organization of the HMP, including the data processing and analytical steps, and the scientific entities gathered to conduct the project. An overview of available HMP data sets and additional resources are provided in Supplementary Tables 1–3. Donors were recruited and enrolled into the HMP through the two clinical centres. Over 240 adults were carefully screened and phenotyped before sampling one to three times at 15 (male) or 18 (female) body sites using a common sampling protocol (<http://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/GetPdf.cgi?id=phd003190.2>). All included subjects were between the ages of 18 and 40 years and had passed a

\*Lists of participants and their affiliations appear at the end of the paper.

screening for systemic health based on oral, cutaneous and body mass exclusion criteria (<http://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/GetPdf.cgi?id=phd002854.2>) (K. Aagaard *et al.*, manuscript submitted).

A Data Analysis and Coordination Center (DACC) was created to serve as the central repository for all HMP WGS, 16S and reference genome sequence information generated by the four sequencing centres. The DACC supports access to analysis software, biological samples, clinical protocols, news, publication announcements and project statistics, and performed centralized analysis of HMP reference genome and WGS annotation in cooperation with the sequencing centres. All unprocessed 16S, WGS and reference genome sequence data are deposited at the National Center for Biotechnology Information (NCBI) (<http://www.ncbi.nlm.nih.gov/bioproject/43021>). Unless otherwise noted, all data sets and protocols described here are available to the scientific community at the DACC (<http://hmpdacc.org>). Specific data sets referred to in this work and available at the DACC are indicated in parentheses with the preface 'RES'.

### Phase I 16S and WGS sequencing overview

A set of 5,298 samples were collected from 242 adults (K. Aagaard *et al.*, manuscript submitted; Table 1 and Supplementary Table 4), from which 16S and WGS data were generated for a total of 5,177 taxonomically characterized communities (16S) and 681 WGS samples describing the microbial communities from habitats within the human airways, skin, oral cavity, gut and vagina. For a subset of 560 samples, both data types were generated (Table 1). These efforts constitute our initial primary metagenomic data sets (phase I) described in more detail later. Additional efforts are ongoing to sequence and analyse the remaining samples from the complete HMP collection (11,174 primary specimens in total from 300 individuals sampled up to three times over 22 months) (K. Aagaard *et al.*, manuscript submitted).

### 16S standards development and sequencing

The goals of the HMP required that 16S sequences and profiles from data produced at the four participating sequencing centres be comparable in a variety of downstream analyses; however, no suitable methodology was available at the commencement of the project. While establishing 16S protocols, we determined that many components of data production and processing can contribute errors and artefacts. We investigated methods that avoid these errors and their subsequent effects on taxonomic classification and operational

taxonomic unit (OTU)-based community structure. The results are discussed in detail in Supplementary Information and ref. 6. Thus, multiple evaluations of 16S protocols were undertaken before adopting a single standardized protocol that ensured consistency in the high-throughput production.

To maximize accuracy and consistency, protocols were evaluated primarily using a synthetic mock community of 21 known organisms<sup>6</sup> (Supplementary Table 5). Additional testing of the protocol was carried out on a subset of HMP samples (Supplementary Table 1). Collectively, these efforts resulted in adoption of a protocol to amplify and sequence samples using the Roche-454 FLX Titanium platform<sup>6</sup> ([http://www.hmpdacc.org/doc/HMP\\_MDG\\_454\\_16S\\_Protocol.pdf](http://www.hmpdacc.org/doc/HMP_MDG_454_16S_Protocol.pdf)). The HMP created both cell mixtures and genomic DNA extracts of the mock community (Supplementary Tables 2 and 5). A large body of metagenomic data (both 16S and WGS) (RES:HMMC) from these and other calibration experiments are available to the community to facilitate further benchmarking of new molecular and analytical approaches (Supplementary Table 3).

The majority of the sample collection was targeted for 16S sequencing using the 454 FLX Titanium based strategy<sup>6</sup>. The nucleotide sequence of the 16S rRNA gene consists of regions of highly conserved sequence, which alternate with nine regions or windows of variable nucleotide sequence that constitute the most informative portions of the gene sequence for use in taxonomic classification. A window covering number three (V3) to five (V5) variable regions (V35) of the 16S rRNA gene was chosen as the target for 4,879 samples. Sequence of a V1 to V3 (V13) window was also included for a subset of 2,971 samples to provide a complementary view of taxonomic profiles<sup>6</sup> (RES:HMR16S) (Table 1, Supplementary Figs 2, 3 and Supplementary Information).

After adoption of the 16S protocol, including removal of multiple sources of potential artefacts or bias generated by 16S sequencing using pyrosequencing<sup>7,8</sup>, a variety of approaches for accurate diversity estimation were developed and compared<sup>9</sup>. A 16S data processing pipeline was established using the mothur software package<sup>10</sup> (Supplementary Information), which includes two optional low and high stringency approaches. The former provides an output favouring longer read lengths tailored towards taxonomic classification, the latter an output with more aggressive sequence error reduction tailored towards OTU construction (RES:HMMCP). A third complementary pipeline was also developed using the QIIME software package<sup>11</sup> (Supplementary Information), which processes these data using an

**Table 1 | HMP donor samples examined by 16S and WGS**

Body region	Body site	Total samples	Total 16S samples	V13 samples	V13 read depth (M)*	V35 samples	V35 read depth (M)*	Samples V13 and V35	Total WGS samples	Total read depth (G)†	Filtered reads (%)‡	Human reads (%)§	Remaining read depth (G)†	Samples 16S and WGS
Gut	Stool	352	337	193	1.4	328	2.4	184	136	1,720.7	15	1	1,450.6	124
Oral cavity	Buccal mucosa	346	330	184	1.3	314	1.7	168	107	1,438.0	9	82	136.7	91
	Hard palate	325	325	179	1.2	310	1.7	164	1	10.9	20	25	5.9	1
	Keratinized gingiva	335	329	183	1.3	319	1.7	173	6	72.3	5	47	34.4	0
	Palatine tonsils	337	332	189	1.2	315	1.9	172	6	74.8	2	80	13.5	1
	Saliva	315	310	166	0.9	292	1.5	148	5	55.7	1	91	4.2	0
	Subgingival plaque	334	328	186	1.2	314	1.8	172	7	92.1	5	79	15.3	1
	Supragingival plaque	345	331	192	1.3	316	1.9	177	115	1,500.7	15	40	674.8	101
	Throat	331	325	176	1.0	312	1.7	163	7	78.8	4	79	13.6	1
	Tongue dorsum	348	332	193	1.3	320	2.0	181	122	1,620.1	15	19	1,084.3	106
	Airway	Anterior nares	316	302	169	1.0	283	1.2	150	84	1,129.9	3	96	14.3
Skin	Left antecubital fossa	269	269	158	0.7	221	0.5	110	0	NA	NA	NA	0	NA
	Left retroauricular crease	313	312	188	1.6	295	1.5	171	9	126.3	9	73	22.1	8
	Right antecubital fossa	274	274	158	0.7	229	0.5	113	0	NA	NA	NA	0	NA
	Right retroauricular crease	319	316	190	1.4	304	1.6	178	15	181.9	18	59	42.4	12
Vagina	Mid-vagina	145	143	91	0.6	140	1.0	88	2	22.6	0	99	0.2	0
	Posterior fornix	152	142	89	0.6	136	1.0	83	53	702.1	6	90	25.2	43
	Vaginal introitus	142	140	87	0.6	131	0.9	78	3	36.5	1	98	0.6	1
Total		5,298	5,177	2,971	19	4,879	26.3	2,673	681	8,863.3	11	49	3,538.1	560

NA, not applicable.

\*  $1 \times 10^6$  reads post-processing with the mothur pipeline (Supplementary Information).

†  $1 \times 10^9$  reads (Supplementary Information).

‡ Fraction of reads with low quality bases that were removed (Supplementary Information).

§ Fraction of human reads that were removed (Supplementary Information).

OTU-binning strategy to which taxonomic classification is added (RES:HMQCP). All pipelines result in highly comparable views of the human microbiome.

### Metagenomic assembly and gene cataloguing

Approximately 749 samples representing targeted body sites were chosen for WGS sequencing using the Illumina GAIIx platform with 101-base-pair paired-end reads. From a high-quality set of 681 samples an average depth of 13 Gb ( $\pm$  4.3) was achieved per sample, collectively producing a total of 8.8 Tb (RES:HMIWGS) (Table 1). Theoretically, these per sample data are sufficient to cover a 3 Mb bacterial genome present at only 0.8% abundance with a probability of 90% (M. C. Wendl *et al.*, manuscript submitted). In addition, 12 stool samples were simultaneously sequenced using the 454 FLX Titanium platform (RES:HM4WGS). Comparisons between the centres demonstrated high consistency of target sequencing depth and success rates<sup>4</sup>. After development of a protocol for removing reads resulting from human DNA contamination (Supplementary Information), 49% of the reads were targeted for removal as human (for information on authorized access to these reads, see Supplementary Information). Samples collected from soft tissue tended to have higher human contamination (for example, mid-vagina (96%), anterior nares (82%) and throat (75%)). Preparations from saliva were also high in human DNA sequence (80%), whereas stool contained a relatively low abundance of human reads (up to 1%) (Supplementary Fig. 4).

After application of a quality control protocol that includes human sequence removal, quality filtering and trimming of reads (Supplementary Information), the remaining 3.5 Tb from 681 samples were subjected to a three-tiered complementary analysis strategy (Supplementary Information) of reference genome mapping (which was able to use  $\sim$ 57% of the data), assembly and gene prediction ( $\sim$ 50% of the data), and metabolic reconstruction ( $\sim$ 36% of the data). This combined strategy facilitated the extraction of maximal organismal and functional information.

Metagenomic assemblies were generated for all available samples using an optimized SOAPdenovo protocol with parameters designed to produce substrates for downstream analyses such as gene and function prediction, resulting in a total of 41 million contigs (RES:HMASM) (Supplementary Information). Reads that remained unassembled were pooled across individual body sites and re-assembled using the same approach, resulting in an additional 4,200,672 contigs (RES:HMBSA). These body-site-specific assemblies are aimed at reconstructing organisms that represent too small a fraction in any individual sample to assemble but are found among many individuals. For 12 stool samples both Illumina and 454 FLX Titanium data (RES:HM4WGS) were generated, allowing a hybrid assembly approach using Newbler (Supplementary Information) (RES:HMHASM). Overall, the assembly statistics recovered varied substantially depending on body site and community complexity (Supplementary Fig. 5). However, our results indicate that, for the assembly strategy we used, metagenomic assembly quality plateaus at approximately 6 Gb of microbial sequence coverage for a sample possessing a microbial community structure similar to that of stool samples (Supplementary Fig. 6).

A WGS-based perspective of community membership was obtained by aligning the reads to a set of 1,742 finished bacterial, 131 archaeal, 3,683 viral and 326 microeukaryotic reference genomes<sup>12</sup> (RES:HMREFG) (Supplementary Information) representing a broad taxonomic range from each of these four domains. A total of 57.6% of the high-quality microbial reads could be associated with a known genome (ranging from 33–77% for anterior nares and posterior fornix, respectively) (RES:HMSCP). The overwhelming majority of mapped sequences originated from bacteria (99.7%), while the remaining reads mapped to microeukaryotes (0.3%) or archaea (<0.01%) (Supplementary Information).

Two complementary approaches were used to summarize overall function and metabolism of the human microbiome, producing two primary data sets of annotations (RES:HMMRC and RES:HMGI) (Supplementary Information) and additional secondary analyses (RES:HMGS, HMHGI, HMGC and HMGOI) (Supplementary Information) available to the community for further interrogation. The first primary data set of annotations was produced by mapping individual shotgun reads to characterized protein families<sup>13</sup> (RES:HMMRC). The second was produced from functionally annotated gene predictions generated from the metagenomic assemblies (RES:HMGI), which were subsequently grouped according to high-level biological processes and to selected additional processes specific to metabolism and regulation<sup>14</sup> (RES:HMGS) (Supplementary Tables 6, 7 and Supplementary Fig. 7).

### HMP data generation and analysis lessons

A key manner in which the HMP resources will serve to guide future studies of the microbiome is by enabling informed decisions regarding sampling protocols and genomic DNA preparation (K. Aagaard *et al.*, manuscript submitted), sequencing depth (M. C. Wendl *et al.*, manuscript submitted), statistical power (P. S. La Rosa *et al.*, manuscript submitted) and metagenomic data type. As indicated in Table 1, the consortium successfully amplified 16S sequences to our target depth at all 18 body sites, with the fewest sequences recovered consistently from the antecubital fossae. The amount of host human DNA recovered and the finest level of OTU resolution varied for 16S sequences among body sites<sup>6</sup> (Supplementary Figs 3 and 4).

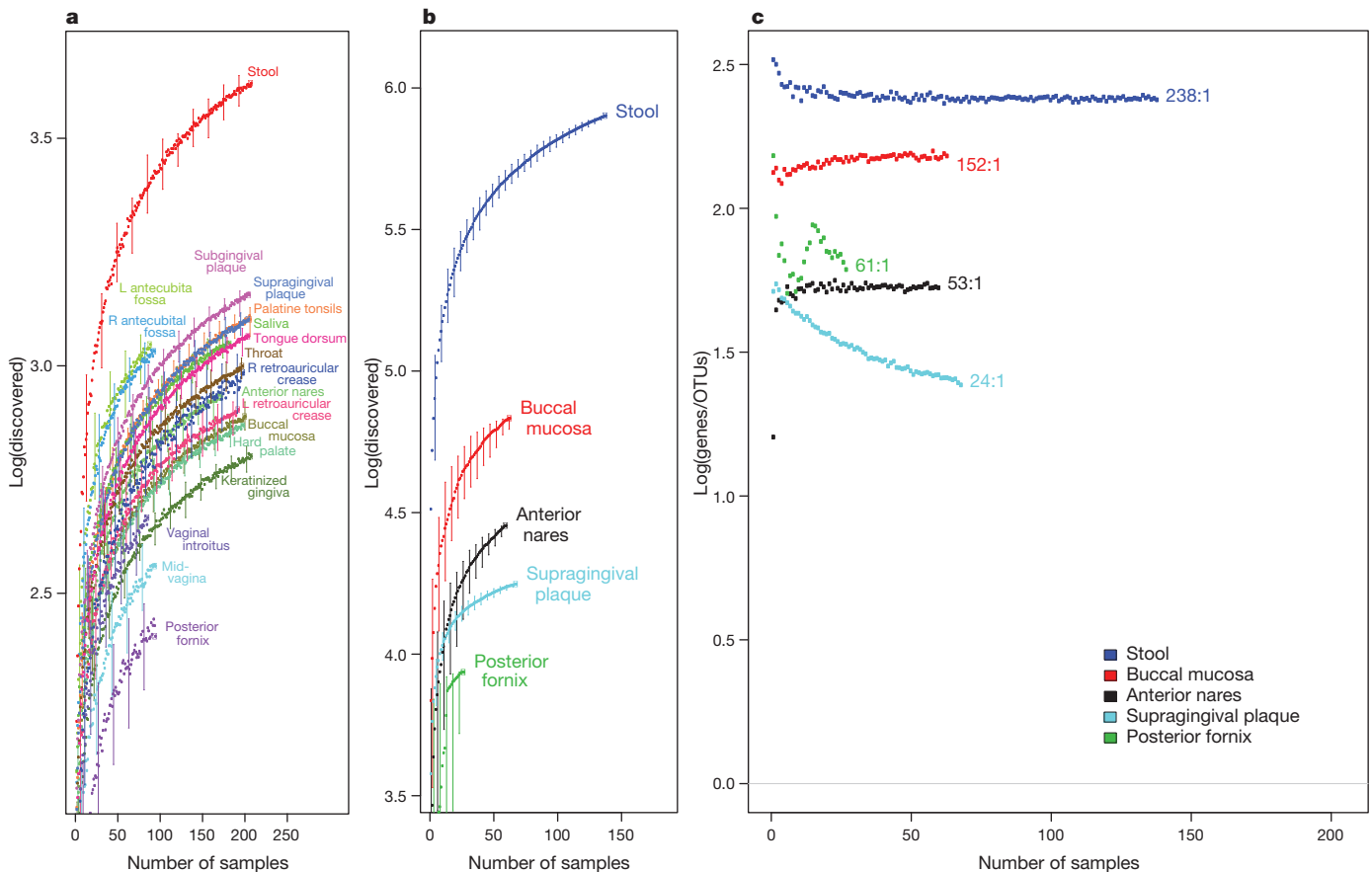
From our WGS investigations, a series of protocols ([http://hmpdacc.org/tools\\_protocols/tools\\_protocols.php](http://hmpdacc.org/tools_protocols/tools_protocols.php)) have been established to process large volumes of short-read WGS data and to annotate and examine these data through both a multi-tiered assembly approach and as single reads<sup>15</sup>. An investigator's choice of metagenomic technologies can thus be guided not only by a 16S versus WGS dichotomy, but also by the expected fraction of host sequence and the appropriate 16S region targeting the dominant taxa at each body site (Supplementary Figs 2–6 and 8).

Together, these data sets represent comprehensive and complementary views of the human microbiome, as shown by comparing organismal (Fig. 1a) and gene (Fig. 1b) catalogues, and the ratio of genes contributed per OTU (Fig. 1c). The discovery rate of new gene clusters (as determined by annotation of assembled WGS data) is in general detected more slowly relative to organismal discovery (as determined by OTU data) owing to the fragmentary nature of these community reads and assemblies despite high sequence depth (Fig. 1a, b and Supplementary Fig. 9), and the number of genes contributed per OTU varies by body site (Fig. 1c and Supplementary Information). However, in general, these results highlight an important point for consideration of further microbiome investigations using these data sets, as they suggest that the majority of the common taxa and genes present in this reference population have been detected.

We additionally compared the gut community gene catalogue sampled by the HMP with that of MetaHIT in terms of total detected gene counts. The HMP recovered more total non-redundant gene counts (5,140,472) than reported by MetaHIT (3,299,822)<sup>5</sup>, probably reflecting a combination of the increased sequence depth obtained by the HMP (11.7 Gb HMP, 4.5 Gb MetaHIT on average) and differences in data generation and processing<sup>5</sup>.

The two non-redundant sets of gene sequences were subsequently combined and compared by matches to a database of orthologous groups<sup>16</sup> of functionally annotated genes. Approximately 57% of the orthologous groups recovered by this method overlapped between the data sets, while an additional 34% versus 10% were unique to the HMP and MetaHIT, respectively (Supplementary Fig. 10, Supplementary Table 8 and Supplementary Information). After removal of genes that received any orthologous group assignment, the remaining novel





**Figure 1 | Rates of gene and OTU discovery from HMP taxonomic and metagenomic data.** a–c, Accumulation curves for OTU counts from 16S data (all body sites) (a), clustered gene index counts from metagenomic data (all applicable body sites) (b) and the ratio of average unique genes contributed versus unique OTUs encountered with increasing sample counts (c) (Supplementary Information). L, left; R, right. Ratios given for each curve in

genes were subsequently clustered<sup>17</sup>. Approximately 79% of the HMP-derived novel gene clusters were orthogonal to one or more clusters in MetaHIT, while an additional 16% were unique to this study versus 5% for MetaHIT-derived data<sup>5</sup> (Supplementary Fig. 11, Supplementary Table 8 and Supplementary Information). These results suggest that, for this body habitat, relatively similar gene catalogues were recovered despite differences in experimental design and protocols. However, a greater proportion of both annotated and unique novel genes were detected in the HMP data set, emphasizing the utility of sequencing depth in recovering gene function and, in particular, deriving rare function. These results further underscore the importance of large-scale sequence-based studies of the microbiome to characterize better its gene content and diversity.

### Human microbiome reference genomes

The current goal for the reference genome component of the HMP is to sequence at least 3,000 reference bacterial genomes, and additional viral and microeukaryotic genomes, associated with the human body. Thus far, more than 800 genomes have been sequenced and are available from the NCBI and the DACC (<http://hmpdacc.org/HMRGD>). From an alignment of WGS reads to reference genomes (RES:HMREFG), approximately 26% from the total read set (46% of all reads that could be aligned) were matched to a subset of 223 HMP reference genomes (Supplementary Information and Supplementary Data).

We continue to solicit community feedback for strains that will best benefit our attempts at understanding the breadth of human microbiome diversity. For example, a prioritized list of the

c represent the average number of unique genes contributed per unique OTU at the final sample count. Curves for stool, buccal mucosa and anterior nares suggest that the proportion of gene-to-taxa discovery has stabilized. In contrast, the curve for supragingival plaque suggests that relatively fewer new genes are being contributed per additional OTU. Error bars represent 95% confidence intervals.

‘most wanted’ HMP taxa is being maintained ([http://hmpdacc.org/most\\_wanted/](http://hmpdacc.org/most_wanted/)) with the goal of targeting these difficult to obtain organisms using both culture-based and single-cell approaches.

A catalogue of all HMP reference genomes along with custom filtering, viewing, graphing and download options can be found at the DACC Project Catalogue ([http://www.hmpdacc-resources.org/hmp\\_catalog/main.cgi](http://www.hmpdacc-resources.org/hmp_catalog/main.cgi)). In addition, comparative analyses of reference genomes are provided by the data warehouse and analytical systems, Integrated Microbial Genomes/HMP ([http://www.hmpdacc-resources.org/cgi-bin/imgm\\_hmp/main.cgi](http://www.hmpdacc-resources.org/cgi-bin/imgm_hmp/main.cgi)). Cultures of all HMP reference strains are required to be made publicly available through the Biodefense and Emerging Infections Research Resources Repository (BEI). Information on strain acquisition can be found at the DACC ([http://hmpdacc.org/reference\\_genomes/reference\\_genomes.php](http://hmpdacc.org/reference_genomes/reference_genomes.php)) and BEI (<http://www.beiresources.org/tabid/1901/stabid/1901/CollectionLinkID/4/Default.aspx>).

### Conclusion

An overarching goal of this multi-year, multi-centre project is the generation of a community resource to advance research efforts related to the microbiome. The result is a collection of 11,174 primary biological specimens representing the human microbiome, as well as corresponding blood samples from the human donors, which are being reserved for sequencing at a future date and from which cell lines will be developed. A variety of new protocols were developed to enable a project of this scope; these include methods for donor recruitment, laboratory and sequence processing, and analysis of

16S and WGS sequence and profiles. These resources serve as models to guide the design of similar projects. Studies with a primary focus on disease can use this reference for comparative purposes, including detecting shifts in microbial taxonomic and functional profiles, or identification of new species not present in healthy cohorts that appear under disease conditions. The catalogue described in this study is, to our knowledge, the largest and most comprehensive reference set of human microbiome data associated with healthy adult individuals. Collectively the data represent a treasure trove that can be mined to identify new organisms, gene functions, and metabolic and regulatory networks, as well as correlations between microbial community structure and health and disease<sup>4</sup>. Among other future benefits, this resource may promote the development of novel prophylactic strategies such as the application of prebiotics and probiotics to foster human health.

## METHODS SUMMARY

As part of a multi-institutional collaboration, the HMP human subjects study was reviewed by the Institutional Review Boards (IRBs) at each sampling site: the BCM (IRB protocols H-22895 (IRB no. 00001021) and H-22035 (IRB no. 00002649)); Washington University School of Medicine (IRB protocol HMP-07-001 (IRB no. 201105198)); and St Louis University (IRB no. 15778). The study was also reviewed by the J. Craig Venter Institute under IRB protocol 2008-084 (IRB no. 00003721), and at the Broad Institute of MIT and Harvard the study was determined to be exempt from IRB review. All study participants gave their written informed consent before sampling and the study was conducted using the Human Microbiome Project Core Sampling Protocol A. Each IRB has a federal-wide assurance and follows the regulations established in 45 CFR Part 46. The study was conducted in accordance with the ethical principles expressed in the Declaration of Helsinki and the requirements of applicable federal regulations.

All further details are in Supplementary Information.

Received 2 November 2011; accepted 10 May 2012.

- Gilbert, J. A. & Dupont, C. L. Microbial metagenomics: beyond the genome. *Annu. Rev. Mar. Sci.* **3**, 347–371 (2011).
- NIH HMP Working Group et al. The NIH Human Microbiome Project. *Genome Res.* **19**, 2317–2323 (2009).
- Human Microbiome Jumpstart Reference Strains Consortium. A catalog of reference genomes from the human microbiome. *Science* **328**, 994–999 (2010).
- The Human Microbiome Project Consortium. Structure, function and diversity of the healthy human microbiome. *Nature* <http://dx.doi.org/10.1038/nature11234> (this issue).
- Qin, J. et al. A human gut microbial gene catalogue established by metagenomic sequencing. *Nature* **464**, 59–65 (2010).
- Jumpstart Consortium Human Microbiome Project Data Generation Working Group. Evaluation of 16S rDNA-based community profiling for human microbiome research. *PLoS ONE* <http://dx.plos.org/10.1371/journal.pone.0039315> (14 June 2012).
- Kunin, V., Engelbrekton, A., Ochman, H. & Hugenholtz, P. Wrinkles in the rare biosphere: pyrosequencing errors can lead to artificial inflation of diversity estimates. *Environ. Microbiol.* **12**, 118–123 (2010).
- Huse, S. M., Huber, J. A., Morrison, H. G., Sogin, M. L. & Welch, D. M. Accuracy and quality of massively parallel DNA pyrosequencing. *Genome Biol.* **8**, R143 (2007).
- Schloss, P. D., Gevers, D. & Westcott, S. L. Reducing the effects of PCR amplification and sequencing artifacts on 16S rRNA-based studies. *PLoS ONE* **6**, e27310 (2011).
- Schloss, P. D. et al. Introducing mothur: open-source, platform-independent, community-supported software for describing and comparing microbial communities. *Appl. Environ. Microbiol.* **75**, 7537–7541 (2009).
- Caporaso, J. G. et al. QIIME allows analysis of high-throughput community sequencing data. *Nature Methods* **7**, 335–336 (2010).
- Martin, J. S. et al. Optimizing read mapping to reference genomes to determine composition and species prevalence in microbial communities. *PLoS ONE* <http://dx.doi.org/10.1371/journal.pone.0036427> (14 June 2012).
- Abubucker, S. et al. Metabolic reconstruction for metagenomic data and its application to the human microbiome. *PLoS Comput. Biol.* <http://dx.doi.org/10.1371/journal.pcbi.1002358> (14 June 2012).
- Ashburner, M. et al. Gene ontology: tool for the unification of biology. *Nature Genet.* **25**, 25–29 (2000).
- Goll, J. et al. A case study for large-scale human microbiome analysis using JCVI's Metagenomics Reports (METAREP). *PLoS ONE* <http://dx.doi.org/10.1371/journal.pone.002904> (14 June 2012).
- Muller, J. et al. eggNOG v2.0: extending the evolutionary genealogy of genes with enhanced non-supervised orthologous groups, species and functional annotations. *Nucleic Acids Res.* **38**, D190–D195 (2010).
- Edgar, R. C. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics* **26**, 2460–2461 (2010).

Supplementary Information is linked to the online version of the paper at [www.nature.com/nature](http://www.nature.com/nature).

**Acknowledgements** The consortium would like to thank our external scientific advisory board: R. Blumberg, J. Davies, R. Holt, P. Ossorio, F. Ouellette, G. Schoolnik and A. Williamson. We would also like to thank our collaborators throughout the International Human Microbiome Consortium, particularly the investigators of the MetaHIT project, for advancing human microbiome research. Data repository management was provided by the NCBI and the Intramural Research Program of the NIH National Library of Medicine. We especially appreciate the generous participation of the individuals from the St Louis, Missouri, and Houston, Texas areas who made this study possible. This research was supported in part by NIH grants U54HG004969 to B.W.B.; U54HG003273 to R.A.G.; U54HG004973 to R.A.G., S.K.H. and J.F.P.; U54HG003067 to E. S. Lander.; U54AI084844 to K.E.N.; N01AI30071 to R. L. Strausberg; U54HG004968 to G.M.W.; U01HG004866 to O.W.; U54HG003079 to R.K.W.; R01HG005969 to C.H.; R01HG004872 to R.K.; R01HG004885 to M.P.; R01HG005975 to P.D.S.; R01HG004908 to Y.Y.; R01HG004900 to M. K. Cho and P. Sankar; R01HG005171 to D.E.H.; R01HG004853 to A.L.M.; R01HG004856 to R.R.; R01HG004877 to R.R.S. and R.M.F.; R01HG005172 to P. Spicer; R01HG004857 to M.P.; R01HG004906 to T.M.S.; R21HG005811 to E.A.-V.; G.A.B. was supported by UH2AI083263 and UH3AI083263 (G.A.B., C. N. Cornelissen, L. K. Eaves and J. F. Strauss); M.J.B. was supported by UH2AR057506, S.M.H. was supported by UH3DK083993 (V. B. Young, E. B. Chang, F. Meyer, T.M.S., M. L. Sogin, J. M. Tiedje); K.P.R. was supported by UH2DK083990 (J.V.); J.A.S. and H.H.K. were supported by UH2AR057504 and UH3AR057504 (J.A.S.); DP2OD001500 to K.M.A.; N01HG62088 to the Coriell Institute for Medical Research; U01DE016937 to F.E.D.; S.K.-H. was supported by RC1DE020298 and R01DE021574 (S.K.-H. and H. Li); J.I. was supported by R21CA139193 (J.I. and D. S. Michaud); K.P.L. was supported by P30DE020751 (D. J. Smith); Army Research Office grant W911NF-11-1-0473 to C.H.; National Science Foundation grants NSF DBI-1053486 to C.H. and NSF IIS-0812111 to M.P.; The Office of Science of the US Department of Energy under contract no. DE-AC02-05CH11231 for P.S.C.; LANL Laboratory-Directed Research and Development grant 20100034DR and the US Defense Threat Reduction Agency grants B1041531 and B0845311 to P.S.C.; Research Foundation - Flanders (FWO) grant to K.F. and J. Raes; R.K. is a Howard Hughes Medical Institute (HHMI) Early Career Scientist; Gordon & Betty Moore Foundation funding and institutional funding from the J. David Gladstone Institutes to K.S.P.; A.M.S. was supported by fellowships provided by the Rackham Graduate School and the NIH Molecular Mechanisms in Microbial Pathogenesis Training Grant T32AI007528; a Crohn's and Colitis Foundation of Canada Grant in Aid of Research to E.A.-V.; 2010 IBM Faculty Award to K.C.W. Analysis of the HMP data was performed using National Energy Research Scientific Computing resources; the BluBioU Computational Resource at Rice University.

**Author Contributions** Principal investigators: B.W.B., R.A.G., S.K.H., B.A.M., K.E.N., J.F.P., G.M.W., O.W., R.K.W. Manuscript preparation: B.A.M., K.E.N., M.P., H.H.C., M.G.G., D.G., C.H., J.F.P. Funding agency management: C.C.B., T.B., V.R.B., J.L.C., S.C., C.D., V.D.F., C.G., M.Y.G., R.D.L., J.M., P.M., J.P., L.M.P., J.A.S., L.W., C.W., K.A.W. Project leadership: S.A., J.H.B., B.W.B., A.T.C., H.H.C., A.M.E., M.G.F., R.S.F., D.G., M.G.G., K.H., S.K.H., C.H., E.A.L., R.M., V.M., J.C.M., B.A.M., M.M., D.M.M., K.E.N., J.F.P., E.J.S., J.V., G.M.W., O.W., A.M.W., K.C.W., J.R.W., S.K.Y., Q.Z. Analysis preparation for manuscript: M.B., B.L.C., D.G., M.G.G., M.E.H., C.H., K.L., B.A.M., X.Q., J.R.W., M.T. Data release: L.A., T.B., I.A.C., K.C., H.C., N.J.D., D.J.D., A.M.E., V.M.F., L.F., J.M.G., S.G., S.K.H., M.E.H., C.J., V.J., C.K., A.A.M., V.M.M., T.M., M.M., D.M.M., J.O., K.P., J.F.P., C.P., X.Q., R.K.S., N.S., I.S., E.J.S., D.V.W., O.W., K.W., K.C.W., C.Y., B.P.Y., Q.Z. Methods and research development: S.A., H.M.A., M.B., D.M.C., A.M.E., R.L.E., M.F., S.F., M.G.F., D.C.F., D.G., G.G., B.J.H., S.K.H., M.E.H., W.A.K., N.L., K.L., V.M., E.R.M., B.A.M., M.M., D.M.M., C.N., J.F.P., M.E.P., X.Q., M.C.R., C.R., E.J.S., S.M.S., D.G.T., D.V.W., G.M.W., Y.W., K.M.W., S.Y., B.P.Y., S.K.Y., Q.Z. DNA sequence production: S.A., E.A., T.A., T.B., C.J.B., D.A.B., K.D.D., S.P.D., A.M.E., R.L.E., C.N.F., S.F., C.C.F., L.L.F., R.S.F., B.H., S.K.H., M.E.H., V.J., C.L.K., S.L.L., N.L., L.L., D.M.M., I.N., C.N., M.O., J.F.P., X.Q., J.G.R., Y.R., M.C.R., D.V.W., Y.W., B.P.Y., Y.Z. Clinical sample collection: K.M.A., M.A.C., W.M.D., L.L.F., N.G., H.A.H., E.L.H., J.A.K., W.A.K., T.M., A.L.M., P.M., S.M.P., J.F.P., G.A.S., J.V., M.A.W., G.M.W. Body site experts: K.M.A., E.A.V., G.A., L.B., M.J.B., C.C.D., F.E.D., L.F., J.I., J.A.K., S.K.H., H.H.K., K.P.L., P.J.M., J. Ravel, T.M.S., J.A.S., J.D.S., J.V. Ethical, legal and social implications: R.M.F., D.E.H., W.A.K., N.B.K., C.M.L., A.L.M., R.R., P. Sankar, P. Spicer, R.R.S., L.Z. Strain management: E.A.V., J.H.B., I.A.C., K.C., S.W.C., H.H.C., T.Z.D., A.S.D., A.M.E., M.G.F., M.G.G., S.K.H., V.J., N.C.K., S.L.L., L.L., K.L., E.A.L., V.M.M., B.A.M., D.M.M., K.E.N., I.N., I.P., L.S., E.J.S., C.M.T., M.T., D.V.W., G.M.W., A.M.W., Y.W., K.M.W., B.P.Y., L.Z. 16S data analysis: K.M.A., E.J.A., G.L.A., C.A.A., M.B., B.W.B., J.P.B., G.A.B., S.R.C., J.C., J.C., T.Z.D., F.E.D., E.D., A.M.E., R.C.E., M.F., A.A.F., J.F., K.F., H.G., D.G., B.J.H., T.A.H., S.M.H., C.H., J.I., J.K.J., S.T.K., S.K.H., R.K., H.H.K., O.K., P.S.L., R.E.L., K.L., C.A.L., D.M., B.A.M., K.A.M., M.M., M.P., J.F.P., M.P., K.S.P., X.Q., J. Raes, K.P.R., M.C.R., B.R., J.F.S., P.D.S., T.M.S., N.S., J.A.S., W.D.S., T.J.S., C.S.S., E.J.S., R.M.T., J.V., T.A.V., Z.W., D.V.W., G.M.W., J.R.W., K.M.W., Y.Y., S.Y., Y.Z. Shotgun data processing and alignments: C.J.B., J.C.C., E.D., D.G., A.G., M.E.H., H.J., D.K., K.C.K., C.L.K., Y.L., J.C.M., B.A.M., M.M., D.M.M., J.O., J.F.P., X.Q., J.G.R., R.K.S., N.U.S., I.S., E.J.S., G.G.S., S.M.S., J.W., Z.W., G.M.W., O.W., K.C.W., T.W., S.K.Y., L.Z. Assembly: H.M.A., C.J.B., P.S.C., L.C., Y.D., S.P.D., M.G.F., M.E.H., H.J., S.K., B.L., Y.L., C.L., J.C.M., J.M.M., J.R.M., P.J.M., M.M., J.F.P., M.P., M.E.P., X.Q., M.R., R.K.S., M.S., D.D.S., G.G.S., S.M.S., C.M.T., T.J.T., W.W., G.M.W., K.C.W., L.Y., Y.Y., S.K.Y., L.Z. Annotation: O.O.A., V.B., C.J.B., I.A.C., A.T.C., K.C., H.H.C., A.S.D., M.G.G., J.M.G., J.G., A.G., S.G., B.J.H., K.H., S.K.H., C.H., H.J., N.C.K., R.M., V.M.M., K.M., T.M., M.M., J.O., K.P., M.P., X.Q., N.S., E.J.S., G.G.S., S.M.S., M.T., G.M.W., K.C.W., J.R.W., C.Y., S.K.Y., Q.Z., L.Z. WGS Metabolic Reconstruction: S.A., B.L.C., J.G., C.H., J.I., B.A.M., M.M., B.R., A.M.S., N.S., M.T., G.M.W., S.Y., Q.Z., J.D.Z.

**Author Information** Accession numbers for all primary sequencing data are given in Supplementary Information. Reprints and permissions information is available at [www.nature.com/reprints](http://www.nature.com/reprints). This paper is distributed under the terms of the Creative Commons Attribution-Non-Commercial-Share-Alike licence, and is freely available to all readers at [www.nature.com/nature](http://www.nature.com/nature). The authors declare no competing financial

interests. Readers are welcome to comment on the online version of this article at [www.nature.com/nature](http://www.nature.com/nature). Correspondence and requests for materials should be addressed to B.A.M. ([bmeth@icvi.org](mailto:bmeth@icvi.org)).

## The Human Microbiome Project Consortium

Barbara A. Methé<sup>1</sup>, Karen E. Nelson<sup>1</sup>, Mihai Pop<sup>2</sup>, Heather H. Creasy<sup>3</sup>, Michelle G. Giglio<sup>3</sup>, Curtis Huttenhower<sup>4,5</sup>, Dirk Gevers<sup>5</sup>, Joseph F. Petrosino<sup>6,15,79</sup>, Sahar Abubucker<sup>7</sup>, Jonathan H. Badger<sup>77</sup>, Asif T. Chinwalla<sup>7</sup>, Ashlee M. Earl<sup>5</sup>, Michael G. FitzGerald<sup>5</sup>, Robert S. Fulton<sup>7</sup>, Kimberlie Hallsworth-Pepin<sup>7</sup>, Elizabeth A. Lobos<sup>7</sup>, Ramana Madupu<sup>1</sup>, Vincent Magrini<sup>7</sup>, John C. Martin<sup>7</sup>, Makedonka Mitreva<sup>7</sup>, Donna M. Muzny<sup>5</sup>, Erica J. Sodergren<sup>7</sup>, James Versalovic<sup>8,9</sup>, Aye M. Wollam<sup>7</sup>, Kim C. Worley<sup>6</sup>, Jennifer R. Wortman<sup>5</sup>, Sarah K. Young<sup>5</sup>, Qiangdong Zeng<sup>5</sup>, Kjersti M. Aagaard<sup>10</sup>, Olukemi O. Abolude<sup>3</sup>, Emma Allen-Vercoc<sup>11</sup>, Eric J. Alm<sup>5,12</sup>, Lucia Alvarado<sup>5</sup>, Gary L. Andersen<sup>13</sup>, Scott Anderson<sup>5</sup>, Elizabeth Appelbaum<sup>7</sup>, Harindra M. Arachchi<sup>5</sup>, Gary Armitage<sup>14</sup>, Cesar A. Arze<sup>3</sup>, Tulin Ayvaz<sup>15</sup>, Carl C. Baker<sup>16</sup>, Lisa Begg<sup>17</sup>, Tsegahiwot Belachew<sup>18</sup>, Veena Bhonegiri<sup>1</sup>, Monika Bihan<sup>1</sup>, Martin J. Blaser<sup>19</sup>, Toby Bloom<sup>5</sup>, Vivien R. Bonazzi<sup>20</sup>, Paul Brooks<sup>21,22</sup>, Gregory A. Buck<sup>22,23</sup>, Christian J. Buhay<sup>6</sup>, Dana A. Busam<sup>1</sup>, Joseph L. Campbell<sup>18,20</sup>, Shane R. Canon<sup>24</sup>, Brandi L. Cantarel<sup>5</sup>, Patrick S. Chain<sup>25,26</sup>, I-Min A. Chen<sup>27</sup>, Lei Chen<sup>27</sup>, Shaila Chhibba<sup>20</sup>, Ken Chu<sup>27</sup>, Dawn M. Ciulla<sup>5</sup>, Jose C. Clemente<sup>28</sup>, Sandra W. Clifton<sup>7</sup>, Sean Conlan<sup>80</sup>, Jonathan Crabtree<sup>3</sup>, Mary A. Cutting<sup>29</sup>, Noam J. Davidovics<sup>3</sup>, Catherine C. Davis<sup>30</sup>, Todd Z. DeSantis<sup>31</sup>, Carolyn Deal<sup>18</sup>, Kimberley D. Delehaunty<sup>7</sup>, Floyd E. Dewhurst<sup>32,33</sup>, Elena Deych<sup>7</sup>, Yan Ding<sup>6</sup>, David J. Dooling<sup>7</sup>, Shannon P. Dugan<sup>6</sup>, W. Michael Dunne Jr<sup>34,35</sup>, A. Scott Durkin<sup>1</sup>, Robert C. Edgar<sup>36</sup>, Rachel L. Erlich<sup>5</sup>, Candace N. Farmer<sup>7</sup>, Ruth M. Farrell<sup>37</sup>, Karoline Faust<sup>38,39</sup>, Michael Feldgarden<sup>5</sup>, Victor M. Felix<sup>3</sup>, Sheila Fisher<sup>5</sup>, Anthony A. Fodor<sup>40</sup>, Larry Forney<sup>41</sup>, Leslie Foster<sup>1</sup>, Valentina Di Francesco<sup>18</sup>, Jonathan Friedman<sup>42</sup>, Dennis C. Friedrich<sup>5</sup>, Catrina C. Fronick<sup>7</sup>, Lucinda L. Fulton<sup>7</sup>, Hongyu Gao<sup>8</sup>, Nathalia Garcia<sup>43</sup>, Georgia Giannoukos<sup>5</sup>, Christina Giblin<sup>18</sup>, Maria Y. Giovanni<sup>18</sup>, Jonathan M. Goldberg<sup>5</sup>, Johannes Goll<sup>1</sup>, Antonio Gonzalez<sup>44</sup>, Allison Griggs<sup>5</sup>, Sharvari Gujja<sup>5</sup>, Brian J. Haas<sup>5</sup>, Holli A. Hamilton<sup>29</sup>, Emily L. Harris<sup>29</sup>, Theresa A. Heppburn<sup>5</sup>, Brandi Herter<sup>7</sup>, Diane E. Hoffmann<sup>45</sup>, Michael E. Holder<sup>6</sup>, Clinton Howarth<sup>5</sup>, Katherine H. Huang<sup>5</sup>, Susan M. Huse<sup>46</sup>, Jacques Izard<sup>32,47</sup>, Janet K. Jansson<sup>48</sup>, Huaiyang Jiang<sup>6</sup>, Catherine Jordan<sup>5</sup>, Vandita Joshi<sup>6</sup>, James A. Katancik<sup>49</sup>, Wendy A. Keitel<sup>15</sup>, Scott T. Kelley<sup>50</sup>, Cristyn Kells<sup>5</sup>, Susan Kinder-Haake<sup>51†</sup>, Nicholas B. King<sup>52</sup>, Rob Knight<sup>28,53</sup>, Dan Knights<sup>44</sup>, Heidi H. Kong<sup>54</sup>, Omry Koren<sup>55</sup>, Sergey Koren<sup>2</sup>, Karthik C. Kota<sup>7</sup>, Christie L. Kovar<sup>6</sup>, Nikos C. Kyrpides<sup>26</sup>, Patricia S. La Rosa<sup>56</sup>, Sandra L. Lee<sup>6</sup>, Katherine P. Lemon<sup>32,57</sup>, Niall Lennon<sup>6</sup>, Cecil M. Lewis<sup>58</sup>, Lora Lewis<sup>6</sup>, Ruth E. Ley<sup>55</sup>, Kelvin Li<sup>1</sup>, Konstantinos Liolios<sup>26</sup>, Bo Liu<sup>2</sup>, Yue Liu<sup>6</sup>, Chien-Chi Lo<sup>25</sup>, Catherine A. Lozupone<sup>28</sup>, R. Dwayne Lunsford<sup>29</sup>, Tessa Madden<sup>59</sup>, Anup A. Mahurkar<sup>5</sup>, Peter J. Mannon<sup>60</sup>, Elaine R. Mardis<sup>7</sup>, Victor M. Markowitz<sup>26,27</sup>, Konstantinos Mavrommatis<sup>26</sup>, Jamison M. McCorriston<sup>1</sup>, Daniel McDonald<sup>28</sup>, Jean McEwen<sup>20</sup>, Amy L. McGuire<sup>61</sup>, Pamela McInnes<sup>29</sup>, Teena Mehta<sup>5</sup>, Kathie A. Mihindukulasuriya<sup>7</sup>, Jason R. Miller<sup>1</sup>, Patrick J. Minx<sup>7</sup>, Irene Newsham<sup>6</sup>, Chad Nusbaum<sup>5</sup>, Michelle O'Laughlin<sup>7</sup>, Joshua Orvis<sup>3</sup>, Ioanna Pagani<sup>26</sup>, Krishna Palaniappan<sup>27</sup>, Shital M. Patel<sup>62</sup>, Matthew Pearson<sup>5</sup>, Jane Peterson<sup>60</sup>, Mircea Podar<sup>63</sup>, Craig Pohl<sup>1</sup>, Katherine S. Pollard<sup>64,65,66</sup>, Margaret E. Priest<sup>5</sup>, Lita M. Proctor<sup>20</sup>, Xiang Qin<sup>6</sup>, Jeroen Raes<sup>38,39</sup>, Jacques Ravel<sup>3,67</sup>, Jeffrey G. Reid<sup>6</sup>, Mina Rho<sup>68</sup>, Rosamond Rhodes<sup>69</sup>, Kevin P. Riehle<sup>70</sup>, Maria C. Rivera<sup>22,23</sup>, Beltran Rodriguez-Mueller<sup>50</sup>, Yu-Hui Rogers<sup>1</sup>, Matthew C. Ross<sup>15</sup>, Carsten Ruder<sup>5</sup>, Ravi K. Sanka<sup>1</sup>, Pamela Sankar<sup>71</sup>, J. Fah Sathirapongasuti<sup>4</sup>, Jeffery A. Schloss<sup>20</sup>, Patrick D. Schloss<sup>72</sup>, Thomas M. Schmidt<sup>73</sup>, Matthew Scholz<sup>25</sup>, Lynn Schriml<sup>3</sup>, Alyxandria M. Schubert<sup>74</sup>, Nicola Segata<sup>4</sup>, Julia A. Segre<sup>80</sup>, William D. Shannon<sup>56</sup>, Richard R. Sharp<sup>37</sup>, Thomas J. Sharpton<sup>64</sup>, Narmada Shenoy<sup>5</sup>, Nihar U. Sheth<sup>22</sup>, Gina A. Simone<sup>74</sup>, Indresh Singh<sup>1</sup>, Chris S. Smillie<sup>42</sup>, Jack D. Sobel<sup>75</sup>, Daniel D. Sommer<sup>2</sup>, Paul Spicer<sup>58</sup>, Granger G. Sutton<sup>1</sup>, Sean M. Sykes<sup>5</sup>, Diana G. Tabbaa<sup>5</sup>, Mathangi Thiagarajan<sup>1</sup>, Chad M. Tomlinson<sup>7</sup>, Manolito Torralba<sup>1</sup>, Todd J. Treangen<sup>76</sup>, Rebecca M. Truty<sup>54</sup>, Tatiana A. Vishnivetskaya<sup>63</sup>, Jason Walker<sup>1</sup>, Lu Wang<sup>20</sup>, Zhengyuan Wang<sup>5</sup>, Doyle V. Ward<sup>5</sup>, Wesley Warren<sup>7</sup>, Mark A. Watson<sup>34</sup>, Christopher Wellington<sup>20</sup>, Kris A. Wetterstrand<sup>20</sup>, James R. White<sup>3</sup>, Katarzyna Wilczek-Boney<sup>6</sup>, Yuan Qing Wu<sup>6</sup>, Kristine M. Wylie<sup>7</sup>, Todd Wylie<sup>7</sup>, Chandri Yandava<sup>5</sup>, Liang Ye<sup>7</sup>, Yuzhen Ye<sup>68</sup>, Shibu Yooseph<sup>77</sup>, Bonnie P. Youmans<sup>15</sup>, Lan Zhang<sup>6</sup>, Yanjiao Zhou<sup>7</sup>, Yiming Zhu<sup>6</sup>, Laurie Zoloth<sup>78</sup>, Jeremy D. Zucker<sup>5</sup>, Bruce W. Birren<sup>5</sup>, Richard A. Gibbs<sup>6</sup>, Sarah K. Highlander<sup>6,15</sup>, George M. Weinstock<sup>7</sup>, Richard K. Wilson<sup>7</sup> & Owen White<sup>3</sup>

<sup>1</sup>J. Craig Venter Institute, 9704 Medical Center Drive, Rockville, Maryland 20850, USA. <sup>2</sup>University of Maryland, Center for Bioinformatics and Computational Biology and Department of Computer Science, Biomolecular Sciences Building Rm. 3120F, College Park, Maryland 20742, USA. <sup>3</sup>University of Maryland School of Medicine, Institute for Genome Sciences 801 W. Baltimore Street, Baltimore, Maryland 21201, USA. <sup>4</sup>Harvard School of Public Health, Department of Biostatistics, 655 Huntington Avenue, Boston, Massachusetts 02115, USA. <sup>5</sup>The Broad Institute of MIT and Harvard, 7 Cambridge Center, Cambridge, Massachusetts 02142, USA. <sup>6</sup>Baylor College of Medicine Human Genome Sequencing Center, One Baylor Plaza, Houston, Texas 77030, USA. <sup>7</sup>Washington University School of Medicine, The Genome Institute, 4444 Forest Park Avenue, St Louis, Missouri 63108, USA. <sup>8</sup>Baylor College of Medicine, Department of Pathology & Immunology, One Baylor Plaza, Houston, Texas 77030, USA. <sup>9</sup>Texas Children's Hospital Department of Pathology, 6621 Fannin Street, Houston, Texas 77030, USA. <sup>10</sup>Baylor College of Medicine, Department of Obstetrics & Gynecology, Division of Maternal-Fetal Medicine, One Baylor Plaza, Houston, Texas 77030, USA. <sup>11</sup>University of Guelph Department of Molecular and Cellular Biology, 50 Stone Road East, Guelph, Ontario N1G 2W1, Canada. <sup>12</sup>Massachusetts Institute of Technology, Department of Civil & Environmental Engineering, Parsons Laboratory, Room 48-317, 15 Vassar Street, Cambridge, Massachusetts 02139, USA. <sup>13</sup>Lawrence Berkeley National Laboratory, Center for Environmental Biotechnology, 1 Cyclotron Road, Berkeley, California 94720, USA. <sup>14</sup>University of California, San Francisco, School of Dentistry, 707 Parnassus Avenue,

San Francisco, California 94143, USA. <sup>15</sup>Baylor College of Medicine, Molecular Virology and Microbiology, One Baylor Plaza, Houston, Texas 77030, USA. <sup>16</sup>National Institutes of Health, National Institute of Arthritis and Musculoskeletal and Skin Diseases (NIAMS), 6701 Democracy Boulevard, MSC 4872, Bethesda, Maryland 20892, USA. <sup>17</sup>National Institutes of Health, Office of Research on Women's Health (ORWH), 6707 Democracy Boulevard, MSC 5484, Bethesda, Maryland 20892, USA. <sup>18</sup>National Institutes of Health, National Institute for Allergy and Infectious Diseases (NIAID), 6610 Rockledge Drive, MSC 6603, Bethesda, Maryland 20892, USA. <sup>19</sup>New York University Langone Medical Center, Department of Medicine, 550 First Avenue, OBV A-606, New York, New York 10016, USA. <sup>20</sup>National Institutes of Health, National Human Genome Research Institute (NHGRI), 5635 Fishers Lane, MSC 9305, Bethesda, Maryland 20892, USA. <sup>21</sup>Virginia Commonwealth University, Department of Statistical Sciences and Operations Research, PO Box 843083, Richmond, Virginia 23284, USA. <sup>22</sup>Virginia Commonwealth University, Center for the Study of Biological Complexity, 1000 West Cary Street, Richmond, Virginia 23284, USA. <sup>23</sup>Virginia Commonwealth University, Department of Biology, 1000 West Cary Street, Richmond, Virginia 23284, USA. <sup>24</sup>Lawrence Berkeley National Laboratory, Technology Integration Group, National Energy Research Scientific Computing Center, 1 Cyclotron Road, Berkeley, California 94720, USA. <sup>25</sup>Los Alamos National Laboratory Genome Science Group, Bioscience Division, HRL, MS-888, LANL, Los Alamos, New Mexico 87545, USA. <sup>26</sup>Joint Genome Institute, 2800 Mitchell Drive, Walnut Creek, California 94598, USA. <sup>27</sup>Lawrence Berkeley National Laboratory, Biological Data Management and Technology Center, Computational Research Division, 1 Cyclotron Road, Berkeley, California 94720, USA. <sup>28</sup>University of Colorado, Department of Chemistry and Biochemistry, Campus Box 215, University of Colorado, Boulder, Colorado 80309-0215, USA. <sup>29</sup>National Institutes of Health, National Institute of Dental and Craniofacial Research (NIDCR), 6701 Democracy Boulevard, MSC 4878, Bethesda, Maryland 20892, USA. <sup>30</sup>The Procter & Gamble Company, FernCare Product Safety and Regulatory Affairs, 6110 Center Hill Avenue, Cincinnati, Ohio 45224, USA. <sup>31</sup>Second Genome, Inc. Bioinformatics Department, 1150 Bayhill Drive, Suite 215, San Bruno, California 94066, USA. <sup>32</sup>Forsyth Institute, Department of Molecular Genetics, 245 First Street, Cambridge, Massachusetts 02142, USA. <sup>33</sup>Harvard School of Dental Medicine, Department of Oral Medicine, Infection and Immunity, 188 Longwood Avenue, Boston, Massachusetts 02115, USA. <sup>34</sup>Washington University School of Medicine, Department of Pathology & Immunology, 660 South Euclid Avenue, Box 8118, St Louis, Missouri 63110, USA. <sup>35</sup>bioMerieux, Inc., 100 Rodolphe Street, Durham, North Carolina 27712, USA. <sup>36</sup>drive5.com, Tiburon, California 94920, USA. <sup>37</sup>Cleveland Clinic, Center for Bioethics, Humanities and Spiritual Care, 9500 Euclid Avenue, Cleveland, Ohio 44195, USA. <sup>38</sup>Vrije Universiteit Brussels, Department of Applied Biological Sciences (DBIT), Pleinlaan 2, 1050 Brussels, Belgium. <sup>39</sup>Vrije Universiteit Brussels, Department of Applied Biological Sciences (DBIT), Pleinlaan 2, 1050 Brussels, Belgium. <sup>40</sup>University of North Carolina Charlotte, Department of Bioinformatics and Genomics, 9201 University City Blvd, Charlotte, North Carolina 28223-0001, USA. <sup>41</sup>University of Idaho, Department of Biological Sciences, Life Sciences South Room 441A, PO Box 443051, Moscow, Idaho 83844, USA. <sup>42</sup>Massachusetts Institute of Technology, Computational and Systems Biology, Parsons Laboratory, Room 48-317, 15 Vassar Street, Cambridge, Massachusetts 02139, USA. <sup>43</sup>Saint Louis University, Center for Advanced Dental Education, 3320 Rutger Street, St Louis, Missouri 63104, USA. <sup>44</sup>University of Colorado, Department of Computer Science, University of Colorado, Boulder, Colorado 80309, USA. <sup>45</sup>University of Maryland Francis King Carey School of Law, 500 W. Baltimore Street, Baltimore, Maryland 21201, USA. <sup>46</sup>Marine Biological Laboratory, Josephine Bay Paul Center, 7 MBL Street, Woods Hole, Massachusetts 02543-1015, USA. <sup>47</sup>Harvard School of Dental Medicine, Department of Oral Medicine, Infection and Immunity, 188 Longwood Avenue, Boston, Massachusetts 02115, USA. <sup>48</sup>Lawrence Berkeley National Laboratory, Ecology Department, Earth Sciences Division, 1 Cyclotron Road, Berkeley, California 94720, USA. <sup>49</sup>University of Texas Health Science Center School of Dentistry, Department of Periodontics, 6516 MD Anderson Blvd, Houston, Texas 77030, USA. <sup>50</sup>San Diego State University, Department of Biology, 5500 Campanile Drive, San Diego, California 92182, USA. <sup>51</sup>UCLA School of Dentistry, Division of Associated Clinical Specialties and Dental Research Institute, 10833 Le Conte Avenue, Los Angeles, California 90095-1668, USA. <sup>52</sup>McGill University, Faculty of Medicine, Peel 3647 Montreal, Quebec H3A 1X1, Canada. <sup>53</sup>Howard Hughes Medical Institute, Campus Box 215, Boulder, Colorado 80309-0215, USA. <sup>54</sup>National Institutes of Health, National Cancer Institute (NCI), Dermatology Branch, CCR, MSC 1908, 10 Center Drive, Bethesda, Maryland 20892, USA. <sup>55</sup>Cornell University, Department of Microbiology, 467 Biotechnology Building, Ithaca, New York 14853, USA. <sup>56</sup>Washington University School of Medicine, Department of Medicine, Division of General Medical Science, 660 South Euclid Avenue, Box 8005, St Louis, Missouri 63110, USA. <sup>57</sup>Children's Hospital Boston, Harvard Medical School, Division of Infectious Diseases, 300 Longwood Avenue, Boston, Massachusetts 02115, USA. <sup>58</sup>University of Oklahoma, Department of Anthropology, 455 West Lindsey, Dale Hall Tower 521, Norman, Oklahoma 73019, USA. <sup>59</sup>Washington University School of Medicine, Department of Obstetrics and Gynecology, 4533 Clayton Avenue, Box 8219, St Louis, Missouri 63110, USA. <sup>60</sup>University of Alabama at Birmingham, Division of Gastroenterology and Hepatology, 1530 3rd Avenue South, Birmingham, Alabama 35294-1150, USA. <sup>61</sup>Baylor College of Medicine, Center for Medical Ethics and Health Policy, One Baylor Plaza, Houston, Texas 77030, USA. <sup>62</sup>Baylor College of Medicine, Medicine-Infectious Disease, One Baylor Plaza, Houston, Texas 77030, USA. <sup>63</sup>Oak Ridge National Laboratory, Biosciences Division, PO Box 2008 MS 6038 Oak Ridge, Tennessee 37831-6038, USA. <sup>64</sup>University of California, San Francisco, Gladstone Institutes, 1650 Owens Street, San Francisco, California 94158, USA. <sup>65</sup>University of California, San Francisco, Institute for Human Genetics, 1650 Owens Street, San Francisco, California 94158, USA. <sup>66</sup>University of California, San Francisco, Division of Biostatistics, 1650 Owens Street, San Francisco, California 94158, USA. <sup>67</sup>University of Maryland School of Medicine, Department of Microbiology and Immunology, BioPark II - Room 611, 801 W. Baltimore Street, Baltimore, Maryland 21201, USA. <sup>68</sup>Indiana University, School of Informatics and Computing, 150 S. Woodlawn Avenue, Bloomington, Indiana 47405, USA. <sup>69</sup>Mount Sinai School of Medicine, Annenberg Building Floor 5th, Room 5-208, 1468 Madison Avenue, New York, New York 10029, USA. <sup>70</sup>Baylor College of Medicine Molecular & Human Genetics, One Baylor Plaza, Houston,



Texas 77030, USA. <sup>71</sup>University of Pennsylvania, Center for Bioethics and Department of Medical Ethics, 3401 Market Street, Suite 320, Philadelphia, Pennsylvania 19104, USA. <sup>72</sup>University of Michigan, Department of Microbiology & Immunology, 5713 Medical Science Bldg. II, 1150 West Medical Center Dr., Ann Arbor, Michigan 48109-5620, USA. <sup>73</sup>Michigan State University, Department of Microbiology and Molecular Genetics, 6180 Biomedical Physical Sciences, Michigan State University, East Lansing, Michigan 48824, USA. <sup>74</sup>The EMMES Corporation, 401 N. Washington St., Suite 700, Rockville, Maryland 20850, USA. <sup>75</sup>Wayne State University School of Medicine, Detroit, Michigan, Harper University Hospital, 3990 John R Street, Detroit, Michigan 48201, USA. <sup>76</sup>Johns Hopkins University School of Medicine, McKusick-Nathans Institute of Genetic Medicine, Bloomberg School of Public Health, E3138, 615 N Wolfe St, Baltimore, Maryland 21205, USA. <sup>77</sup>J. Craig Venter Institute, 10355 Science Center Drive, San Diego, California 92121, USA. <sup>78</sup>Northwestern University, Feinberg School of Medicine, 420 East Superior Street Chicago, Illinois 60611, USA. <sup>79</sup>Alkek Center for Metagenomics and Microbiome Research, Baylor College of Medicine, One Baylor Plaza, Houston, Texas 77030, USA. <sup>80</sup>National Institutes of Health, National Human Genome Research Institute (NHGRI), Genetics and Molecular Biology Branch, MSC 4442, Bethesda, Maryland 20892, USA. †Deceased.