

The status of the archaeal and bacterial census: an update

Running title: The archaeal and bacterial census

Patrick D. Schloss^{1†}, Rene Girard², Thomas Martin², Joshua Edwards², and J. Cameron Thrash^{2†}

† To whom correspondence should be addressed: pschloss@umich.edu and thrashc@lsu.edu

1. Department of Microbiology and Immunology, University of Michigan, Ann Arbor, MI 48109

2. Department of Biological Sciences, Louisiana State University, Baton Rouge, LA 70803

Abstract

A census is typically carried out for people across a range of geographical levels; however, microbial ecologists have implemented a molecular census of bacteria and archaea by sequencing their 16S rRNA genes. We assessed how well the census of full-length 16S rRNA gene sequences is proceeding in the context of recent advances in high throughput sequencing technologies because full-length sequences are typically used as references for classification of the short sequences generated by newer technologies. Among the 1,411,234 and 53,546 full-length bacterial and archaeal sequences, 94.5% and 95.1% of the bacterial and archaeal sequences, respectively, belonged to operational taxonomic units (OTUs) that have been observed more than once. Although these metrics suggest that the census is approaching completion, 29.2% of the bacterial and 38.5% of the archaeal OTUs have been observed more than once. Thus, there is still considerable diversity to be explored. Unfortunately, the rate of new full-length sequences has been declining and new sequences are primarily being deposited by a small number of studies. Furthermore, sequences from soil and aquatic environments, which are known to be rich in bacterial diversity, only represent 7.8 and 16.5% of the census while sequences associated with host-associated environments represent 55.0% of the census. Continued use of traditional approaches and new technologies such as single cell genomics and short read assembly are likely to improve our ability to sample rare OTUs if it is possible to overcome this sampling bias. The success of ongoing efforts to use short read sequencing to characterize archaeal and bacterial communities requires that researchers strive to expand the depth and breadth of this census.

Importance

The biodiversity contained within the bacterial and archaeal domains dwarfs that of the eukaryotes and the services these organisms provide to the biosphere are critical. Surprisingly, we have done a relatively poor job of formally tracking the quality of the biodiversity as represented in full-length 16S rRNA genes. By understanding how this census is proceeding, it is possible to suggest the best allocation of resources for advancing the census. We found that the ongoing effort has done

27 an excellent job of sampling the most abundant organisms, but struggles to sample the more rare
28 organisms. Through the use of new sequencing technologies we should be able to obtain full-length
29 sequences from these rare organisms. Furthermore, we suggest that by allocating more resources
30 to sampling environments known to have the greatest biodiversity we will be able to make significant
31 advances in our characterization of archaeal and bacterial diversity.

Introduction

The effort to quantify the number of different organisms in a system remains fundamental to understanding ecology (1, 2). At the scale of microorganisms, small physical sizes, morphological ambiguity, and highly variable population sizes complicate this process. Furthermore, creating standards for delimiting what makes one microbe “different” from another has been contentious (3, 4). In spite of these challenges, we continue to peel back the curtain on the microbial world with the aid of more and more informative, if still limited, technologies like cultivation, 16S rRNA gene surveys, single cell technologies, and metagenomics.

Generating a comprehensive understanding of any system with a single gene may seem a fool’s errand, yet we have learned a considerable amount regarding the diversity, dynamics, and natural history of microorganisms using the venerable 16S rRNA gene. In 1983, the full-length 16S rRNA gene sequence of *Escherichia coli* (accession J01695) was released as part of NCBI’s GenBank making it the first of what is now more than 10 million 16S rRNA gene sequences to be deposited into the database (5). 16S rRNA gene accessions represent nearly one-third of all sequences deposited in GenBank, making it the best-represented gene. As Sanger sequencing has given way to so-called “next generation sequencing” technologies, hundreds of millions of 16S rRNA gene sequences have been deposited into the NCBI’s Sequence Read Archive. The expansion in sequencing throughput and increased access to sequencing technology has allowed for more environments to be sequenced at a deeper coverage, resulting in the identification of novel taxa. The ability to obtain sequence data from microorganisms without cultivation has radically altered our perspective of their role in nearly every environment from deep ocean sediment cores (e.g. accession AY436526) to the International Space Station (e.g. accession DQ497748).

Previously, Schloss and Handelsman (6) assigned the 56,215 partial 16S rRNA gene sequences that were available in the Ribosomal Database Project to operational taxonomic units (OTUs) and concluded that the sampling methods of the time were insufficient to identify the previously estimated 10^7 to 10^9 different species (7, 8). That census called for a broader and deeper characterization of all environments. Refreshingly, this challenge was largely met. There have been major investments in studying the Earth’s microbiome using 16S rRNA gene sequencing through initiatives such as the

Human Microbiome Project (9), the Earth Microbiome Project (10), and the International Census of Marine Microorganisms (11). But most importantly, the original census was performed on the cusp of radical developments in sequencing technologies. That advancement has moved the generation of sequencing throughput from large sequencing centers to individual investigators and leveraged their diverse interests to expand the representation of organisms and environments represented in public databases.

It is disconcerting that the increase in sequencing volume has come at the cost of sequence length. The commonly used MiSeq-based sequencing platform from Illumina is extensively used to sequence the approximately 250 bp V4 hypervariable region of the 16S rRNA gene; other schemes have used different parts of the gene that are generally shorter than 500 bp. The number of OTUs that are inferred from the sequencing data when using different regions within the 16S rRNA gene can vary considerably and the genetic diversity within these regions typically has only a modest correlation with the genetic diversity of the full-length sequence (12, 13). Thus, it remains unclear to what degree richness estimates from short read technologies over or underestimate the numbers from full-length sequences. Furthermore, we likely lack the full-length reference sequences necessary to adequately classify the novel biodiversity we are sampling when we generate 100-times the sequence data from a community than we did using full-length sequencing.

Here we update the status of the archaeal and bacterial census with full-length 16S rRNA gene sequences. In the 13 years since the collection of data for Schloss and Handelsman's initial census, the number of full-length sequences has grown exponentially, despite the overwhelming contemporary focus by most researchers on short-read technologies. This update to the census allows us to evaluate the relative sampling thoroughness for different environments and clades and make an argument for the continued need to collection full-length sequence data from many systems that have a long history of study. As researchers consider coalescing into a Unified Microbiome Initiative (14), it will be important to balance the need for mechanism-based studies with the need to generate full-length reference sequences from a diversity of environments.

Results and Discussion

The status of the bacterial and archaeal census. To assess the field's progress in characterizing the biodiversity of bacteria and archaea, we assigned each 16S rRNA gene sequence to OTUs using distance thresholds that varied between 0 and 20%. Although it is not possible to link a specific taxonomic level (e.g. species, genus, family, etc.) to a specific distance threshold, we selected distances of 0, 3, 5, 10, and 20% because they are widely regarded as representing the range of genetic diversity of the 16S rRNA gene within each domain. By rarefaction, it was clear that the ongoing sampling efforts have started to saturate the number of current OTUs. After sampling 1,411,234 near full-length bacterial 16S rRNA gene sequences we have identified 217,645, 108,950, 66,819, 15,743, and 3,731 OTUs at the respective thresholds (Figure 1A, Table 1). Using only the OTUs generated using a 3% threshold, we calculated a 94.5% Good's coverage (percent of sequences belonging to OTUs that have been observed more than once), but only 29.2% OTU coverage (percent of the OTUs that have been observed more than once). Paralleling the bacterial results, after sampling 53,546 archaeal 16S rRNA gene sequences we have identified 11,040, 4,252, 2,364, 812, and 110 OTUs (Figure 1B, Table 1). Using only the OTUs generated with a 3% threshold, we calculated a 95.1% Good's coverage, but only 38.5% OTU coverage. These results indicate that regardless of the domain, continued sampling with the current strategies for generating full-length sequences will largely reveal OTUs that have already been observed, even though a large fraction of OTUs have only been sampled once. Considering more than 70.8% of the OTUs have only been observed once, it is likely that an even larger number of OTUs have yet to be sampled for both domains.

Sequencing efforts are a source of bias in the census. One explanation for the large number of OTUs that have only been observed once is that with the the broad adoption of sequencing platforms that generate short sequence reads, the rate of full-length sequence generation has declined. In fact, since 2009 the number of new bacterial sequences generated has slowed to an average of 189,960 sequences per year (Figure 2A). Although this is still an impressive number of sequences, since 2007 the number of new bacterial OTUs has plateaued at an average of 11,184 new OTUs per year (Figure 2B). Given the expense of generating full-length sequences using the

Sanger sequencing technology and the transition to other platforms at that time, we expected that the large number of sequences were being deposited by a handful of large projects. Indeed, when we counted the number of submissions responsible for depositing 50% of the sequences, we found that with the exception of 2006 and 2013, eight or fewer studies were responsible for depositing the majority of the full-length sequences each year since 2005 (Figure 2C). Between 2009 and 2012, 908,190 total sequences were submitted and 6 submissions from 5 studies were responsible for depositing 550,274 (60.6% of all sequences). These studies generated sequences from the human gastrointestinal tract (15), human skin (16, 17), murine skin (18), and hypersaline microbial mats (19). The heavy focus on host-associated communities is reflected in the rarefaction curve for this category (Figure 1C). In contrast to recent years, between 1995 and 2006, an average of 39.3 studies were responsible for submitting more than half of the sequences each year. Although the recent deep surveys represent significant contributions to our knowledge of bacterial biogeography, their small number and lack of environmental diversity is indicative of the broader problems in advancing the bacterial census.

The depth of sequencing being done to advance the archaeal census has been 26-times less than that of the bacterial census (Table 1). The annual number of sequences submitted has largely paralleled that of the bacterial census with a plateau starting in 2009 and an average of 7,075 sequences each year since then. The number of new archaeal OTUs represented by these sequences began to slow in 2005 with an average of 355.5 new OTUs per year. With the exception of 2012 and 2014, the number of submissions responsible for more than 50% of the archaeal sequences submitted per year has varied between 2 and 11 submissions per year. The clear bias towards sequencing bacterial 16S rRNA genes has limited the ability to more fully characterize the biodiversity of the archaea, which is clearly reflected in the relatively meager sampling effort across habitats, compared to bacteria (Figure 1D),

The ability to sample archaea and bacteria is taxonomically skewed. The Firmicutes, Proteobacteria, Actinobacteria, and Bacteroidetes represent 89.2% of the bacterial sequences and the Euryarchaeota and Thaumarchaeota 86.5% of the archaeal sequences. We sought to understand how the representation of individual phyla has changed relative to the state of the census in 2006. We used 2006 as a reference point for calibrating the dynamics of the bacterial

and archaeal censuses since that was the year that the first highly parallelized 16S rRNA gene sequence dataset was published (20). Based on the representation of sequences within the SILVA database, in 2006 there were 61 bacterial and 18 archaeal phyla. Since then there have been 4 new bacterial (CKC4, OC31, S2R-29, and SBYG-2791) and 2 new archaeal candidate phyla (Ancient Archaeal Group and TVG8AR30). Relative to the overall sequencing trends before and after 2006, several phyla stand out for being over and underrepresented in sequence submissions (Figure 3). Among the bacterial phyla with at least 1,000 sequences, Atribacteria and Kazan-3B-09 were sequenced 4-fold more often while *Deinococcus-Thermus* and *Tenericutes* were sequenced 2-fold less often than would have been expected since 2006. Among the archaeal phyla with at least 1,000 sequences, the Thaumarchaeota were sequenced 2.0-fold more often and the Crenarchaeota were sequenced 6.7-fold less often than expected. Together, these results demonstrate a change in the phylum-level lineages represented in the census from before and after 2006 and encouragingly, show that some underrepresented phyla are becoming better sampled.

Focusing the census by environment. We were able to assign 89.3 and 95.1% of the sequences to one of seven broad environmental categories based on the metadata that accompanied the SILVA database (Tables 1). Across these broad categories there was wide variation in the number of sequences that have been sampled. Among bacterial sequences, the three best represented groups were from host-associated (N=804,585), aquatic (N=214,085), and built environment (N=108,799) sources. Among the archaeal sequences the three best represented groups were the same, but ordered differently: aquatic (N=34,400), built environment (N=7,286), and host-associated (N=5,597) (Figure 1C,D)). For both domains, soil samples were the fourth most represented category (bacteria: 74,870; archaea: 2,517). The orders of these categories was surprising considering soil and aquatic environments harbor the most microbial biomass and biodiversity (21). In spite of wide variation in sequencing depth and coverage (Table 1), the interquartile range across the fine-level categories for the bacterial OTU coverage only varied between 34.5 to 40.0 (median coverage=36.7%). The interquartile range in the OTU coverage by environment for the archaeal data was 41.5 to 53.1 (median coverage=44.9%). The archaeal coverage was higher than that of the bacterial OTU coverage for all categories except the food-associated, plant surface, and other invertebrate categories. Across all categories, the bacterial and archaeal sequencing

data represented a limited number of phyla (Figure 4). Among the bacterial data, the fine-scale categories were dominated by Proteobacteria (N=24), Firmicutes (N=2), and Actinobacteria (N=1) and among the archaeal data, they were dominated by Euryarchaeota (N=16), Thaumarchaeota (N=10), and Aenigmarchaeota (N=1). Regardless, there were clear phylum-level signatures that differentiated the various categories. Within each of the bacterial and archaeal phyla, there was considerable variation in the relative abundance of each across the categories confirming that taxonomic signatures exist to differentiate different environments even at a broad taxonomic level.

The cultured census. In the 2004 bacterial census, there was concern expressed that although culture-independent methods were significantly enhancing our knowledge of microbial life, there were numerous bacterial phyla with no or only a few cultured representatives. To update this assessment, we identified those sequences that came from cultured and uncultured organisms. Overall, 18.9% of bacterial sequences and 6.8% of archaeal sequences have come from isolated organisms. Comparing the fraction of sequences deposited during and before 2006 from isolates to those collected after 2006, we found that culturing rates lag by 2.4 and 2.5-fold for bacteria and archaea, respectively. Among the 65 bacterial phyla, 24 have no cultured representatives and 14 of the 20 archaeal phyla have no cultured representatives. This lag is likely due to the differences in throughput of culture-dependent and -independent approaches. Of the phyla with at least one cultured representative, the median percentage of sequences coming from a culture was only 2.8% for the bacterial phyla and 1.7% for the archaeal phyla (Figure 5). Even though many phyla have cultured representatives, there is still a skew in the representation of most phyla found in cultivation efforts.

Considering the possibility that large culture-independent sequencing efforts may only be re-sequencing organisms that already exist in culture, we asked what percentage of OTUs had at least one cultured representative. We found that 16.9% of the 117,385 bacterial OTUs and 13.1% of the 4,574 archaeal OTUs had at least one cultured representative (Figure 5). Comparing the percentage of sequences with cultured representatives to the percentage of OTUs containing a sequence from a cultured representative revealed a strong cultivation bias within the Firmicutes, which had a higher percentage of sequences generated by cultivated representatives than would be expected based on the number of cultured organisms represented by OTUs (Figure 5). This likely

reflects the extremely high number of cultivated biomedically relevant cultivars from genera such as *Bacillus*, *Streptococcus*, *Lactobacillus*, *Staphylococcus*, and others. Conversely, many phyla, including Cyanobacteria, Actinobacteria, Bacteroidetes, and Nitrospirae, had a lower percentage of sequences belonging to cultivated representatives than would be expected based on the percentage of OTUs that have sequences from cultured organisms, indicating that the cultivation efforts in these clades are relatively inefficient with regards to available diversity. Nevertheless, it is clear that the majority of OTUs from any phylum remain uncultivated, to say nothing of the diversity of organisms that may be encapsulated within the 97% sequence identity cutoff.

New technologies to access novel biodiversity. Given the shift from Sanger sequencing to platforms that offer higher throughput but shorter reads, we are concerned that our ability to harvest full-length sequences from communities will remain stalled. Several culture-independent methods have been developed that offer the ability to obtain full-length sequences of the 16S rRNA gene and even the complete genome. These have included single cell genomics (22) and assembly of short 16S rRNA gene fragments using data generated from PCR amplicons or metagenomic shotgun sequence data with the Expectation-Maximization Iterative Reconstruction of Genes from the Environment (EMIRGE) algorithm (23, 24). To test the ability of these technologies to expand our knowledge of microbial diversity beyond that of traditional approaches, we compared the overlap of OTUs found using each of the new methods with the traditional approaches (Figure 6). Utilizing the 16S rRNA gene sequences extracted from the single-cell genomes available on the Integrated Microbial Genomes (IMG) system (25), we identified 311 bacterial and 70 archaeal sequences, which were assigned to 115 and 27 bacterial and archaeal OTUs, respectively. Interestingly, only 8.7 and 3.7% of the bacterial and archaeal single celled OTUs, respectively, had not been observed by previous efforts. Next, we identified six studies that utilized EMIRGE to assemble 16S rRNA gene sequences from metagenomic sequences (23, 26–30). Together these studies assembled 599 bacterial and 9 archaeal full-length sequences, which were assigned to 335 and 7 bacterial and archaeal OTUs, respectively. Only 40.6 and 60.3% of the bacterial OTUs generated by this approach were previously identified by this traditional cultivation and PCR-based approaches, respectively. Although the application of this approach to Archaea has been limited, it was still surprising that 85.7 and 85.7% of the archaeal OTUs had been previously recovered by traditional

cultivation and PCR-based approaches, respectively. Finally, we pooled 76,080 bacterial sequences from five studies that utilized EMIRGE to assemble 16S rRNA gene sequences from fragmented amplicons (24, 31–34). These sequences were assigned to 40,213 OTUs. We were surprised that only 7.6% of these OTUs were previously found by a more traditional approach. Although these PCR-based EMIRGE results may be valid, the high degree of novelty that was observed suggests that the error of the assembled reads may be too high for generating reference sequences. Each of these methods represent promising opportunities to continue the bacterial census using full-length sequences as well as genomic information.

Conclusions

It is clear that considerable biodiversity has been discovered since the first census in 2004. However, much of it has been biased towards particular phyla and environments. Our analysis suggests that 94.5% of new full-length bacterial and archaeal sequences are likely to have already been seen. Meanwhile, 29.2% of bacterial and 38.5% of archaeal OTUs have only been observed once. In spite of current estimates suggesting the global bacterial species richness may be as high as 10^{12} species (35), the current census based on full-length 16S rRNA gene sequences suggests that existing sampling methods will prevent us from acquiring full-length sequences for that level of diversity. As we have shown, current strategies repeatedly sample the same OTUs and do a poor job of resampling rarer populations. Given this low level of OTU coverage, it is likely that there are many more bacterial and archaeal populations yet to be sampled.

There are several additional reasons to suspect that the current census should be considered conservative. First, we found that most sequences recently deposited into public databases are being made by a small number of projects that have deeply sampled similar environments, and the number of full-length reads deposited into the databases has stalled. Second, it is widely acknowledged that 16S rRNA gene primers are biased; these biases are amplified when designing primers to amplify subregions used in sequencing short reads (36). Assembly of metagenomic data has shown the presence of introns in the 16S rRNA genes of organisms within the so-called “Candidate Phyla Radiation” (e.g. Saccharibacteria (TM7), Peregrinibacteria,

Berkelbacteria (ACD58), WWE3 Microgenomates (OP11), Parcubacteria (OD1), et al.) that would preclude detection with standard PCR-based approaches (37, 38). Third, the willingness of researchers to contribute their sequences and the metadata describing the environment that the sequences were sampled from is critical for assessing the progress of the census and to accrue the benefits from having full-length sequences in the databases. As an illustration of this problem, only 5 of the 11 studies that used the EMIRGE algorithm deposited their sequences in GenBank. This makes the sequences from the other studies effectively invisible to the search algorithms used by 16S rRNA gene-specific databases to harvest sequences. As assembly and long read technologies advance, a mechanism is needed to assess the quality of the consensus sequences and to make them easily accessible to the 16S rRNA gene-specific databases.

Efforts to census archaea and bacteria using short read technology such as the International Census of Marine Microbes, the Earth Microbiome Project, and the Human Microbiome Project have significantly advanced our knowledge of archaeal and bacterial biogeography; however, these analyses have demonstrated the limitations of databases and taxonomies that are based on sequences from common and abundant organisms. During the period prior to the introduction of massively parallelized high throughput sequencing, it was common for a study to generate dozens or hundreds of sequences per sample. The existing databases that are used for classifying sequences are based on those sequences, which represent organisms that are generally abundant. We hypothesize that recent difficulties obtaining adequate classification for short sequences captured from more rare organisms are because our databases do not contain full-length references for those sequences. We fear that these trends will worsen unless researchers can leverage new sequencing and cultivation technologies to generate large numbers of full-length sequences from a large number of diverse samples.

Novel technologies such as single-cell genomics, metagenomics, and algorithms to recover full-length sequences from new sequencing platforms have demonstrated promise in circumventing previous limitations in identifying new OTUs. Using EMIRGE to assemble fragmented 16S rRNA gene amplicons may allow us to obtain deep coverage of communities; however, it is still unclear how faithful the assembled sequence is to that of the original organism. Additional sequencing technologies also offer the ability to directly generate full-length sequences, such as PacBio and

potentially Oxford Nanopore. Initial application of PacBio to sequencing full-length fragments suggests that the sequences suffered from a high error rate (39). To obtain a more direct investigation of rare organisms, microbiologists are developing novel cultivation and single cell genomics techniques (???, 40–42). The ability to enrich or select for specific populations using these approaches could limit the need for redundant brute force sequencing. These approaches are still in active development, and we hope that through continuous refinement, they may allow us to significantly improve the coverage of OTUs in public databases.

Materials and Methods

Sequence data curation. The July 19, 2015 release of the ARB-formatted SILVA small subunit (SSU) reference database (SSU Ref v.123) was downloaded from http://www.arb-silva.de/fileadmin/silva_databases/release_123/ARB_files/SSURef_123_SILVA_19_07_15_opt.arb.tgz (43). This release is based on the EMBL-EBI/ENA Release 123, which was released in March 2015. The SILVA curators identify potential SSU sequences using keyword searches and sequence-based search using RNAmmer (<http://www.arb-silva.de/documentation/release-123/>). The SILVA curators then screened the 7,168,241 resulting sequences based on a minimum length criteria (<300 nt), number of ambiguous base calls (>2%), length of sequence homopolymers (>2%), presence of vector contamination (>2%), low alignment quality value (<75), and likelihood of being chimeric (Pintail value < 50). Of the remaining sequences, the bacterial reference set retained those bacterial sequences longer than 1,200 nucleotides and the archaeal reference set retained those archaeal sequences longer than 900 nucleotides. The aligned 1,515,024 bacterial and 59,240 archaeal sequences were exported from the database using ARB along with the complete set of metadata. Additional sequence data was included from single-cell genomes available on the Integrated Microbial Genomes (IMG) system (25), many of which were recently obtained via the GEBA-MDM effort in Rinke et al. (22). “SCGC” was searched on the IMG database March 12, 2015 to download the bacterial (N=249) and archaeal (N=46) 16S rRNA gene sequences and their associated metadata. Further, sequences generated from amplicon and shotgun metagenomic data using the EMIRGE program were also included (23, 24). The IMG and EMIRGE sequences

were aligned against the respective SILVA-based reference using mothur (44). The aligned bacterial and archaeal sequence sets were pooled and processed in parallel. Using mothur, sequences were further screened to remove any sequence with more than 2 ambiguous base calls and trimmed to overlap the same alignment coordinates. The sequences in the resulting bacterial dataset overlapped bases 113 through 1350 of an *E. coli* reference sequence (V00348) and had a median length of 1,233 nt. The sequences in the resulting archaeal dataset overlapped positions 362 to 937 of a *Sulfolobus solfataricus* reference sequence (X03235) and had a median length of 580 nt. The archaeal sequences were considerably shorter than their initial length because it was necessary to find a common overlapping region across the sequences. The final datasets contained 1,411,234 bacterial and 53,546 archaeal 16S rRNA gene sequences. Sequences were assigned to OTUs using the average neighbor clustering algorithm (45).

Metadata curation. The metadata that was contained within the SSU Ref database was used to expand our analysis beyond a basic count of sequences and the number of OTUs in each domain. The environmental origins of the 16S rRNA gene sequences were manually classified using seven broad “coarse” categories, and further refined to facilitate additional analyses with twenty-six more specific “fine” categories (Table S1). These were assigned based on manual curation of the “isolation_source” category within the ARB database associated with each of the sequences. For source definitions that were not identifiable by online searches, educated guesses were made or they were placed into the coarse “Other” category. There were 151,669 bacterial and 2,565 archaeal sequences where an “isolation_source” term was not collected. We ascertained whether a sequence came from a cultured organism by including those sequences that had data in their “strain” or “isolate” fields within the database and excluded any sequences that had “Unc” as part of their database name as this is a convention in the database that represents sequences from uncultured organisms. Complete tables containing the ARB-provided metadata, taxonomic information, OTU assignment, and our environmental categorizations are available at FigShare for the bacterial (<https://dx.doi.org/10.6084/m9.figshare.2064927>) and archaeal (<https://dx.doi.org/10.6084/m9.figshare.2064942>) data.

Calculating coverage. Sequencing coverage (C_{Sequence}) was quantified by two methods. The first was to use Good’s coverage according to

$$C_{Sequence} = 1 - \frac{n_1}{N_t}$$

where n_1 is the number of OTUs represented by only one sequence and N_t is the total number of sequences (46). Although Good's coverage provides information about the success of the sequencing effort in sampling the most abundant organisms in a community, it does not directly provide information about the success of the sequencing effort in recovering previously unobserved OTUs. To quantify the ability of sequencing to identifying novel OTUs or, in other words, to quantify the "distance" in the peak of the rarefaction curves to their hypothetical asymptote, we defined "OTU coverage" (C_{OTU}) as

$$C_{OTU} = 1 - \frac{n_1}{S_t}$$

where S_t is the total number of OTUs. Whereas Good's coverage estimates the probability that a new sequence will have already been seen, OTU coverage estimates the probability that a new OTU will match an existing one. It is therefore an extension of Good's coverage in that it quantifies the probability that, for any given set of sequences clustered into an OTU, that OTU will have already been seen. Thus, high Good's coverage means that any new sequence is unlikely to be novel, and high OTU coverage means that any new OTU is unlikely to be novel.

Data analysis. Our analysis made use of ARB (OS X v.6.0) (43), mothur (v.1.37.0) (44), and R (v.3.2.0) (47). Within R we utilized the knitr (v.1.10.5), wesanderson (v.0.3.2), and openxlsx (v.2.4.0) packages. A reproducible version of this manuscript including data extraction and processing is available at https://www.github.com/SchlossLab/Schloss_Census2_mBio_2016.

Figure 1. Number of OTUs sampled among bacterial and archaeal 16S rRNA gene sequences for different OTU definitions and level of sequencing effort. Rarefaction curves for different OTU definitions of Bacteria (A) and Archaea (B). Rarefaction curves for the coarse environments in Table 1 for Bacteria (C) and Archaea (D).

Figure 2. Progression of the archaeal and bacterial census since the first full-length 16S rRNA gene sequence was deposited into GenBank in 1983.* The number of bacterial and archaeal 16S rRNA gene sequences deposited (A) and the new OTUs they represent (B) has increased exponentially until the last several years when the rate of change has plateaued. For both bacterial and archaeal sequences, the number of studies that are responsible for depositing more than 50% of the sequences each year has been relatively small (C).

Figure 3. Relative rate of sequence deposition for each bacterial and archaeal phylum before and after 2006 relative to the sequencing of all bacteria. The figure shows the relative rates for those phyla with at least 1,000 sequences and the x-axis is on a log2 scale. The data for all bacterial and archaeal phyla are available in Supplemental Tables 2 and 3, respectively.

Figure 4. Heatmap depicting the relative abundance of the most common bacterial and archaeal phyla across different environments. Each environmental category exhibited a phylum-level signature although the bacterial census was dominated by sequences from the Firmicutes, Proteobacteria, Actinobacteria, and Bacteroidetes and the archaeal census was dominated by sequences from the Euryarchaeota and Thaumarchaeota. The ten most abundant phyla across all environmental categories are shown. The data for all bacterial and archaeal phyla are available in Supplemental Tables 4 and 5, respectively.

Figure 5. The rate that sequences and OTUs are generated from bacterial and archaeal cultures relative to all sequences and OTUs by phylum. Phyla with greater than 1,000 sequences are listed by domain. Open circles indicate the percentage of sequences in the database that match cultured organisms. Closed circles indicate the percentage of OTUs in this analysis that contain sequences belonging to a cultured organism. The data for all bacterial and archaeal phyla are available in Supplemental Tables 6 and 7, respectively.

Figure 6. The percentage of bacterial and archaeal OTUs found by single cell genomics and EMIRGE using PCR or metagenomics that were also detected by other methods. The bars comparing a method to itself indicate the percentage of OTUs that were only detected by that method.

390 **Supplemental Table 1. Description of environmental categories and the criteria used to**
391 **assign sequences to each category.**

392 **Supplementary Table 2. Frequency that each bacterial phylum was sequenced before and**
393 **after 2006.**

394 **Supplementary Table 3. Frequency that each archaeal phylum was sequenced before and**
395 **after 2006.**

396 **Supplementary Table 4. Frequency that each bacterial phylum was found across each of**
397 **the environmental categories.**

398 **Supplementary Table 5. Frequency that each archaeal phylum was found across each of**
399 **the environmental categories.**

400 **Supplementary Table 6. Frequency that each bacterial sequence or OTU was retrieved by**
401 **cultivation or by culture-independent methods.**

402 **Supplementary Table 7. Frequency that each archaeal sequence or OTU was retrieved by**
403 **cultivation or by culture-independent methods.**

References

1. **McGill BJ, Etienne RS, Gray JS, Alonso D, Anderson MJ, Benecha HK, Dornelas M, Enquist BJ, Green JL, He F, Hurlbert AH, Magurran AE, Marquet PA, Maurer BA, Ostling A, Soykan CU, Ugland KI, White EP.** 2007. Species abundance distributions: Moving beyond single prediction theories to integration within an ecological framework. *Ecology Letters* **10**:995–1015. doi:<http://doi.org/10.1111/j.1461-0248.2007.01094.x>.
2. **Hubbell SP.** 2001. *A Unified Theory of Biodiversity and Biogeography*. Princeton University Press, Princeton.
3. **Konstantinidis KT, Ramette A, Tiedje JM.** 2006. The bacterial species definition in the genomic era. *Philosophical Transactions of the Royal Society B: Biological Sciences* **361**:1929–1940. doi:<http://doi.org/10.1098/rstb.2006.1920>.
4. **Oren A, Garrity GM.** 2013. Then and now: A systematic review of the systematics of prokaryotes in the last 80 years. *Antonie van Leeuwenhoek* **106**:43–56. doi:<http://doi.org/10.1007/s10482-013-0084-1>.
5. **Brosius J, Palmer ML, Kennedy PJ, Noller HF.** 1978. Complete nucleotide sequence of a 16S ribosomal RNA gene from *Escherichia coli*. *Proceedings of the National Academy of Sciences* **75**:4801–4805. doi:<http://doi.org/10.1073/pnas.75.10.4801>.
6. **Schloss PD, Handelsman J.** 2004. Status of the microbial census. *Microbiology and Molecular Biology Reviews* **68**:686–691. doi:<http://doi.org/10.1128/mnbr.68.4.686-691.2004>.
7. **Dykhuizen DE.** 1998. Santa Rosalia revisited: Why are there so many species of bacteria? *Antonie van Leeuwenhoek* **73**:25–33. doi:<http://doi.org/10.1023/a:1000665216662>.
8. **Curtis TP, Sloan WT, Scannell JW.** 2002. Estimating prokaryotic diversity and its limits. *Proceedings of the National Academy of Sciences* **99**:10494–10499. doi:<http://doi.org/10.1073/pnas.142680199>.
9. **The Human Microbiome Consortium.** 2012. Structure, function and diversity of the healthy

human microbiome. *Nature* **486**:207–214. doi:<http://doi.org/10.1038/nature11234>.

10. **Gilbert JA, Jansson JK, Knight R.** 2014. The Earth Microbiome Project: Successes and aspirations. *BMC Biology* **12**:69. doi:<http://doi.org/10.1186/s12915-014-0069-1>.

11. **Amaral-Zettler L, Artigas LF, Baross J, P.A. LB, Boetius A, Chandramohan D, Herndl G, Kogure K, Neal P, Pedrós-Alió C, Ramette A, Schouten S, Stal L, Thessen A, Leeuw J de, Sogin M.** 2010. A global census of marine microbes, pp. 221–245. *In* *Life in the worlds oceans*. Wiley-Blackwell.

12. **Youssef N, Sheik CS, Krumholz LR, Najjar FZ, Roe BA, Elshahed MS.** 2009. Comparison of species richness estimates obtained using nearly complete fragments and simulated pyrosequencing-generated fragments in 16S rRNA gene-based environmental surveys. *Applied and Environmental Microbiology* **75**:5227–5236. doi:<http://doi.org/10.1128/aem.00592-09>.

13. **Schloss PD.** 2010. The effects of alignment quality, distance calculation method, sequence filtering, and region on the analysis of 16S rRNA gene-based studies. *PLoS Comput Biol* **6**:e1000844. doi:<http://doi.org/10.1371/journal.pcbi.1000844>.

14. **Alivisatos AP, Blaser MJ, Brodie EL, Chun M, Dangl JL, Donohue TJ, Dorrestein PC, Gilbert JA, Green JL, Jansson JK, Knight R, Maxon ME, McFall-Ngai MJ, Miller JF, Pollard KS, Ruby EG, Taha SA.** 2015. A unified initiative to harness Earth's microbiomes. *Science* **350**:507–508. doi:<http://doi.org/10.1126/science.aac8480>.

15. **Li E, Hamm CM, Gulati AS, Sartor RB, Chen H, Wu X, Zhang T, Rohlf FJ, Zhu W, Gu C, Robertson CE, Pace NR, Boedeker EC, Harpaz N, Yuan J, Weinstock GM, Sodergren E, Frank DN.** 2012. Inflammatory bowel diseases phenotype, *textitC. difficile* and NOD2 genotype are associated with shifts in human ileum associated microbial composition. *PLoS ONE* **7**:e26284. doi:<http://doi.org/10.1371/journal.pone.0026284>.

16. **Kong HH, Oh J, Deming C, Conlan S, Grice EA, Beatson MA, Nomicos E, Polley EC, Komarow HD, Murray PR, Turner ML, Segre JA.** 2012. Temporal shifts in the skin microbiome associated with disease flares and treatment in children with atopic dermatitis. *Genome Research*

22:850–859. doi:<http://doi.org/10.1101/gr.131029.111>.

17. **Grice EA, Kong HH, Conlan S, Deming CB, Davis J, Young AC, Bouffard GG, Blakesley RW, Murray PR, Green ED, Turner ML, Segre JA.** 2009. Topographical and temporal diversity of the human skin microbiome. *Science* **324**:1190–1192. doi:<http://doi.org/10.1126/science.1171700>.

18. **Grice EA, Snitkin ES, Yockey LJ, Bermudez DM, Liechty KW, Segre JA, Mullikin J, Blakesley R, Young A, Chu G, Ramsahoye C, Lovett S, Han J, Legaspi R, Fuksenko T, Reddix-Dugue N, Sison C, Gregory M, Montemayor C, Gestole M, Hargrove A, Johnson T, Myrick J, Riebow N, Schmidt B, Novotny B, Gupti J, Benjamin B, Brooks S, Coleman H, Ho S-I, Schandler K, Smith L, Stantripop M, Maduro Q, Bouffard G, Dekhtyar M, Guan X, Masiello C, Maskeri B, McDowell J, Park M, Thomas PJ.** 2010. Longitudinal shift in diabetic wound microbiota correlates with prolonged skin defense response. *Proceedings of the National Academy of Sciences* **107**:14799–14804. doi:<http://doi.org/10.1073/pnas.1004204107>.

19. **Harris JK, Caporaso JG, Walker JJ, Spear JR, Gold NJ, Robertson CE, Hugenholtz P, Goodrich J, McDonald D, Knights D, Marshall P, Tufo H, Knight R, Pace NR.** 2012. Phylogenetic stratigraphy in the guerrero negro hypersaline microbial mat. *The ISME Journal* **7**:50–60. doi:<http://doi.org/10.1038/ismej.2012.79>.

20. **Sogin ML, Morrison HG, Huber JA, Welch DM, Huse SM, Neal PR, Arrieta JM, Herndl GJ.** 2006. Microbial diversity in the deep sea and the underexplored “rare biosphere”. *Proceedings of the National Academy of Sciences* **103**:12115–12120. doi:<http://doi.org/10.1073/pnas.0605127103>.

21. **Whitman WB, Coleman DC, Wiebe WJ.** 1998. Prokaryotes: The unseen majority. *Proceedings of the National Academy of Sciences* **95**:6578–6583. doi:<http://doi.org/10.1073/pnas.95.12.6578>.

22. **Rinke C, Schwientek P, Sczyrba A, Ivanova NN, Anderson IJ, Cheng J-F, Darling A, Malfatti S, Swan BK, Gies EA, Dodsworth JA, Hedlund BP, Tsiamis G, Sievert SM, Liu W-T, Eisen JA, Hallam SJ, Kyrpides NC, Stepanauskas R, Rubin EM, Hugenholtz P, Woyke T.** 2013. Insights into the phylogeny and coding potential of microbial dark matter. *Nature* **499**:431–437. doi:<http://doi.org/10.1038/nature12352>.

- 482 23. **Miller CS, Baker BJ, Thomas BC, Singer SW, Banfield JF.** 2011. EMIRGE: Reconstruction
483 of full-length ribosomal genes from microbial community short read sequencing data. *Genome Biol*
484 **12**:R44. doi:<http://doi.org/10.1186/gb-2011-12-5-r44>.
- 485 24. **Miller CS, Handley KM, Wrighton KC, Frischkorn KR, Thomas BC, Banfield JF.** 2013.
486 Short-read assembly of full-length 16S amplicons reveals bacterial diversity in subsurface sediments.
487 *PLoS ONE* **8**:e56018. doi:<http://doi.org/10.1371/journal.pone.0056018>.
- 488 25. **Markowitz VM, Chen I-MA, Palaniappan K, Chu K, Szeto E, Pillay M, Ratner A, Huang**
489 **J, Woyke T, Huntemann M, Anderson I, Billis K, Varghese N, Mavromatis K, Pati A, Ivanova**
490 **NN, Kyrpides NC.** 2013. IMG 4 version of the integrated microbial genomes comparative analysis
491 system. *Nucleic Acids Research* **42**:D560–D567. doi:<http://doi.org/10.1093/nar/gkt963>.
- 492 26. **Wrighton KC, Thomas BC, Sharon I, Miller CS, Castelle CJ, VerBerkmoes NC, Wilkins**
493 **MJ, Hettich RL, Lipton MS, Williams KH, Long PE, Banfield JF.** 2012. Fermentation, hydrogen,
494 and sulfur metabolism in multiple uncultivated bacterial phyla. *Science* **337**:1661–1665. doi:<http://doi.org/10.1126/science.1224041>.
495
- 496 27. **Rocha UN da, Cadillo-Quiroz H, Karaoz U, Rajeev L, Klitgord N, Dunn S, Truong V,**
497 **Buenrostro M, Bowen BP, Garcia-Pichel F, Mukhopadhyay A, Northen TR, Brodie EL.** 2015.
498 Isolation of a significant fraction of non-phototroph diversity from a desert biological soil crust. *Front*
499 *Microbiol* **6**:277. doi:<http://doi.org/10.3389/fmicb.2015.00277>.
- 500 28. **Hamilton TL, Jones DS, Schaperdoth I, Macalady JL.** 2015. Metagenomic insights into
501 S(0) precipitation in a terrestrial subsurface lithoautotrophic ecosystem. *Front Microbiol* **5**:756.
502 doi:<http://doi.org/10.3389/fmicb.2014.00756>.
- 503 29. **Handley KM, VerBerkmoes NC, Steefel CI, Williams KH, Sharon I, Miller CS, Frischkorn**
504 **KR, Chourey K, Thomas BC, Shah MB, Long PE, Hettich RL, Banfield JF.** 2012. Biostimulation
505 induces syntrophic interactions that impact c, s and n cycling in a sediment microbial community.
506 *The ISME Journal* **7**:800–816. doi:<http://doi.org/10.1038/ismej.2012.148>.
- 507 30. **Gladden JM, Allgaier M, Miller CS, Hazen TC, VanderGheynst JS, Hugenholtz P,**

- 508 **Simmons BA, Singer SW.** 2011. Glycoside hydrolase activities of thermophilic bacterial
509 consortia adapted to switchgrass. *Applied and Environmental Microbiology* **77**:5804–5812.
510 doi:<http://doi.org/10.1128/aem.00032-11>.
- 511 **31. Brooks B, Firek BA, Miller CS, Sharon I, Thomas BC, Baker R, Morowitz MJ, Banfield JF.**
512 2014. Microbes in the neonatal intensive care unit resemble those found in the gut of premature
513 infants. *Microbiome* **2**:1. doi:<http://doi.org/10.1186/2049-2618-2-1>.
- 514 **32. Wilkins MJ, Wrighton KC, Nicora CD, Williams KH, McCue LA, Handley KM, Miller CS,**
515 **Giloteaux L, Montgomery AP, Lovley DR, Banfield JF, Long PE, Lipton MS.** 2013. Fluctuations
516 in species-level protein expression occur during element and nutrient cycling in the subsurface.
517 *PLoS ONE* **8**:e57819. doi:<http://doi.org/10.1371/journal.pone.0057819>.
- 518 **33. Handley KM, Wrighton KC, Miller CS, Wilkins MJ, Kantor RS, Thomas BC, Williams KH,**
519 **Gilbert JA, Long PE, Banfield JF.** 2014. Disturbed subsurface microbial communities follow
520 equivalent trajectories despite different structural starting points. *Environ Microbiol* **17**:622–636.
521 doi:<http://doi.org/10.1111/1462-2920.12467>.
- 522 **34. Alessi DS, Lezama-Pacheco JS, Janot N, Suvorova EI, Cerrato JM, Giammar DE,**
523 **Davis JA, Fox PM, Williams KH, Long PE, Handley KM, Bernier-Latmani R, Bargar JR.**
524 2014. Speciation and reactivity of uranium products formed during in situ bioremediation
525 in a shallow alluvial aquifer. *Environmental Science & Technology* **48**:12842–12850.
526 doi:<http://doi.org/10.1021/es502701u>.
- 527 **35. Locey KJ, Lennon JT.** 2015. Scaling laws predict global microbial diversity. *PeerJ PrePrints*.
528 doi:<http://doi.org/10.7287/peerj.preprints.1451v1>.
- 529 **36. Parada AE, Needham DM, Fuhrman JA.** 2015. Every base matters: Assessing small subunit
530 rRNA primers for marine microbiomes with mock communities, time series and global field samples.
531 *Environ Microbiol* n/a–n/a. doi:<http://doi.org/10.1111/1462-2920.13023>.
- 532 **37. Brown CT, Hug LA, Thomas BC, Sharon I, Castelle CJ, Singh A, Wilkins MJ, Wrighton**
533 **KC, Williams KH, Banfield JF.** 2015. Unusual biology across a group comprising more than 15%

of domain bacteria. *Nature* **523**:208–211. doi:<http://doi.org/10.1038/nature14486>.

38. **Eloe-Fadrosh EA, Ivanova NN, Woyke T, Kyrpides NC.** 2016. Metagenomics uncovers gaps in amplicon-based detection of microbial diversity. *Nat Microbiol* **1**:5032. doi:<http://doi.org/10.1038/nmicrobiol.2015.32>.

39. **Schloss PD, Westcott SL, Jenior ML, Highlander SK.** 2015. Sequencing 16S rRNA gene fragments using the pacBio sMRT DNA sequencing system. *PeerJ PrePrints*. doi:<http://doi.org/10.7287/peerj.preprints.778v1>.

40. **Nichols D, Cahoon N, Trakhtenberg EM, Pham L, Mehta A, Belanger A, Kanigan T, Lewis K, Epstein SS.** 2010. Use of iChip for high-throughput in situ cultivation of “uncultivable” microbial species. *Applied and Environmental Microbiology* **76**:2445–2450. doi:<http://doi.org/10.1128/aem.01754-09>.

41. **Buerger S, Spoering A, Gavrish E, Leslin C, Ling L, Epstein SS.** 2012. Microbial scout hypothesis, stochastic exit from dormancy, and the nature of slow growers. *Applied and Environmental Microbiology* **78**:3221–3228. doi:<http://doi.org/10.1128/aem.07307-11>.

42. **Das N, Tripathi N, Basu S, Bose C, Maitra S, Khurana S.** 2015. Progress in the development of gelling agents for improved culturability of microorganisms. *Front Microbiol* **6**:698. doi:<http://doi.org/10.3389/fmicb.2015.00698>.

43. **Pruesse E, Quast C, Knittel K, Fuchs BM, Ludwig W, Peplies J, Glockner FO.** 2007. SILVA: A comprehensive online resource for quality checked and aligned ribosomal RNA sequence data compatible with ARB. *Nucleic Acids Research* **35**:7188–7196. doi:<http://doi.org/10.1093/nar/gkm864>.

44. **Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski RA, Oakley BB, Parks DH, Robinson CJ, Sahl JW, Stres B, Thallinger GG, Horn DJV, Weber CF.** 2009. Introducing mothur: Open-source, platform-independent, community-supported software for describing and comparing microbial communities. *Applied and Environmental Microbiology* **75**:7537–7541. doi:<http://doi.org/10.1128/aem.01541-09>.

- 560 45. **Westcott SL, Schloss PD.** 2015. De novo clustering methods outperform reference-based
561 methods for assigning 16S rRNA gene sequences to operational taxonomic units. PeerJ **3**:e1487.
562 doi:<http://doi.org/10.7717/peerj.1487>.
- 563 46. **Good IJ.** 1953. The population frequencies of species and the estimation of population
564 parameters. Biometrika **40**:237–264. doi:<http://doi.org/10.1093/biomet/40.3-4.237>.
- 565 47. **R Core Team.** 2015. R: A language and environment for statistical computing. R Foundation
566 for Statistical Computing, Vienna, Austria.

Table 1. Status of microbial census by habitat classifications and domain. The isolation_source field from the SILVA reference database was manually curated to assign bacterial and archaeal sequences coarse and fine scale habitat classifications. We calculated the number of sequences and OTUs observed and the percent coverage on a sequence or OTU basis for each classification and domain. Descriptions of each category are provided in Table S1.

Coarse	Fine	Bacteria				Archaea			
		Sequences (N)	OTUs (N)	% Seq. Coverage	% OTU Coverage	Sequences (N)	OTUs (N)	% Seq. Coverage	% OTU Coverage
Aerosol		3,472	1,068	79.5	33.2	2	1	100.0	100.0
Aquatic	Brackish	1,094	646	54.6	23.1	1,368	314	87.4	44.9
	Brackish sediment	390	243	54.4	26.7	525	208	76.8	41.3
	Freshwater	21,647	6,689	80.8	37.7	1,540	439	84.7	46.5
	Freshwater sediment	6,733	3,549	63.0	29.8	1,324	488	79.3	43.9
	Marine	134,727	14,287	94.3	46.7	10,983	830	95.8	44.5
	Marine sediment	27,801	9,567	79.6	40.8	14,049	1,507	95.0	53.7
	Hydrothermal vent	10,860	4,216	75.4	36.5	3,797	734	90.4	50.3
	Ice	2,073	936	71.2	36.1	42	5	95.2	60.0
	Other	8,760	3,802	71.7	34.8	772	313	80.7	52.4
Built	Digesters	33,152	8,949	82.9	36.8	4,764	483	93.6	36.4
	Food-associated	11,813	1,632	92.0	41.9	117	40	80.3	42.5
	Industrial/mining	16,582	6,099	76.6	36.3	1,245	336	84.4	42.3
	Pollution associated	38,696	10,602	84.1	41.9	716	249	79.2	40.2
	Other	8,556	2,730	79.1	34.7	444	111	90.8	63.1
Plant associated	Root	19,695	5,052	84.3	38.7	200	61	85.5	52.5
	Surface	4,892	1,385	82.7	38.8	0	0	NA	NA
	Other	9,753	3,217	78.8	35.8	22	7	90.9	71.4
Soil	Agriculture	10,051	4,017	73.6	34.0	146	56	80.8	50.0
	Desert	3,042	1,280	73.7	37.5	245	79	77.6	30.4
	Permafrost	1,922	870	73.0	40.3	39	20	64.1	30.0
	Other	59,855	17,166	82.9	40.4	2,087	516	89.1	55.8
Host-associated	Vertebrate	773,045	42,497	96.1	29.5	5,389	454	95.1	41.6
	Arthropod	13,209	3,688	81.8	34.7	87	52	58.6	30.8
	Other invertebrate	7,476	2,626	78.0	37.3	67	30	73.1	40.0
	Other	10,855	1,754	89.2	33.4	54	17	87.0	58.8
Other		19,414	5,930	81.6	39.9	882	249	84.2	44.2
No source data		151,669	14,144	94.9	45.6	2,565	559	88.6	47.6
Total		1,411,234	108,950	94.5	29.2	53,546	4,252	95.1	38.5