

RESEARCH

Open Access



Species-level bacterial community profiling of the healthy sinonasal microbiome using Pacific Biosciences sequencing of full-length 16S rRNA genes

Joshua P. Earl^{1†}, Nithin D. Adappa^{2†}, Jaroslaw Krol^{1†}, Archana S. Bhat¹, Sergey Balashov¹, Rachel L. Ehrlich¹, James N. Palmer², Alan D. Workman², Mariel Blasetti², Bhaswati Sen¹, Jocelyn Hammond¹, Noam A. Cohen², Garth D. Ehrlich^{1*} and Joshua Chang Mell^{1*}

Abstract

Background: Pan-bacterial 16S rRNA microbiome surveys performed with massively parallel DNA sequencing technologies have transformed community microbiological studies. Current 16S profiling methods, however, fail to provide sufficient taxonomic resolution and accuracy to adequately perform species-level associative studies for specific conditions. This is due to the amplification and sequencing of only short 16S rRNA gene regions, typically providing for only family- or genus-level taxonomy. Moreover, sequencing errors often inflate the number of taxa present. Pacific Biosciences' (PacBio's) long-read technology in particular suffers from high error rates per base. Herein, we present a microbiome analysis pipeline that takes advantage of PacBio circular consensus sequencing (CCS) technology to sequence and error correct full-length bacterial 16S rRNA genes, which provides high-fidelity species-level microbiome data.

Results: Analysis of a mock community with 20 bacterial species demonstrated 100% specificity and sensitivity with regard to taxonomic classification. Examination of a 250-plus species mock community demonstrated correct species-level classification of > 90% of taxa, and relative abundances were accurately captured. The majority of the remaining taxa were demonstrated to be multiply, incorrectly, or incompletely classified. Using this methodology, we examined the microgeographic variation present among the microbiomes of six sinonasal sites, by both swab and biopsy, from the anterior nasal cavity to the sphenoid sinus from 12 subjects undergoing trans-sphenoidal hypophysectomy. We found greater variation among subjects than among sites within a subject, although significant within-individual differences were also observed. *Propionibacterium acnes* (recently renamed *Cutibacterium acnes*) was the predominant species throughout, but was found at distinct relative abundances by site.

(Continued on next page)

* Correspondence: ge33@drexel.edu; jcm385@drexel.edu

[†]Joshua P Earl, Nithin D Adappa and Jaroslaw Krol contributed equally to this work.

¹Department of Microbiology & Immunology, Centers for Genomic Sciences and Advanced Microbial Processing, Drexel University College of Medicine, 245 N 15th Street, Philadelphia, PA 19102, USA

Full list of author information is available at the end of the article



(Continued from previous page)

Conclusions: Our microbial composition analysis pipeline for single-molecule real-time 16S rRNA gene sequencing (MCSMRT, <https://github.com/jpearl01/mcsmrt>) overcomes deficits of standard marker gene-based microbiome analyses by using CCS of entire 16S rRNA genes to provide increased taxonomic and phylogenetic resolution. Extensions of this approach to other marker genes could help refine taxonomic assignments of microbial species and improve reference databases, as well as strengthen the specificity of associations between microbial communities and dysbiotic states.

Keywords: Microbiome, 16S rRNA, Paranasal sinuses, Sinonasal, Database, Long-read DNA sequencing, Circular consensus sequencing

Background

The advent of culture- and cloning-free methods to analyze bacterial phylogenetic marker genes by deep sequencing ushered in a new era of microbial community analysis, dramatically reducing the labor and cost of profiling the identities and abundances of microbes from different environments, independent of their ability to be cultivated [1–4]. The small subunit ribosomal RNA gene (16S rRNA) is shared by all bacteria and has been sequenced in thousands of distinct named species. Because of this, polymerase chain reactions (PCR) using primers that target conserved regions can amplify variable segments of the 16S rRNA gene from across the bacterial domain for amplicon-based deep sequencing [5, 6]. 16S sequence databases can then be used to classify a given sequence read's taxonomic source. Combined with increasingly powerful ecological methods for analyzing microbial community dynamics and inferring community-level metabolic networks, profiling the taxonomic composition of bacterial communities by 16S rRNA gene sequencing has become a standard part of microbiome analysis [7–10].

Unfortunately, the use of the 16S rRNA gene as a taxonomic marker has, in part, been constrained by the short read length of the most commonly used sequencing platform for microbial community profiling (the Illumina MiSeq), which only allows interrogation of up to three of nine variable regions in the 16S rRNA gene (called V1–V9), often targeting only V3–V5, V1–V3, or V4 alone [9, 11–15]. This constraint limits the taxonomic resolution to which short reads can be classified, typically only to the family- or genus-level, and furthermore, taxonomic resolution varies for different groups of bacteria when using different portions of the 16S rRNA gene [9]. Low-resolution classification in turn limits not only the accuracy and precision of ecological inferences and metabolic reconstructions, but also the ability to identify appropriate bacterial strains to use in follow-up experimental and translational studies. Metagenomic shotgun sequencing has been shown to often provide high taxonomic and phylogenetic resolution [16, 17], but these approaches continue to be prohibitively expensive in many cases (particularly when in the presence of excess host

DNA), and consensus remains in flux regarding the best pipelines for shotgun metagenomics-based community analysis [16].

An alternative is to use “3rd generation” long-read sequencing technology to obtain full-length 16S rRNA gene sequences (V1–V9, hereafter FL16S). This increases taxonomic and phylogenetic resolution by increasing the number of informative sites sequenced, while continuing to use a well-studied pan-bacterial marker gene. Initial applications of Pacific Biosciences (PacBio) single-molecule real-time (SMRT) sequencing were hampered by the technology's high intrinsic error rate [18–20], but improvements to the chemistry have since allowed for the generation of high-quality “circular consensus sequence” (CCS) reads, in which individual 16S rRNA genes are sequenced many times using circularized library templates combined with highly processive polymerases that provide for single-molecule, consensus-sequence error correction [21]. Recent studies evaluating FL16S sequencing by PacBio have found that, with appropriate processing and filtering, CCS reads of FL16S genes can be generated that are of sufficiently high quality to offer higher taxonomic resolution than partial 16S rRNA sequences [22–25].

The composition of the human sinonasal microbiome and how it changes in health and disease remains poorly understood, largely due to differences in methodology among studies resulting in large variations in reported bacterial profiles [26–31]. Culture-based approaches capture < 15% of resident bacterial taxa when compared to nucleic acid-based techniques, since fast-growing bacteria like staphylococci tend to predominate in culture specimens, and recovery of anaerobes and slow-growing bacteria is limited [28, 32, 33]. Comparing across recent surveys of the sinonasal bacterial community reveals broadly similar results, but few specific assertions can be made; agreement between studies and results has been limited by an inability to distinguish bacteria at the species level [34–38] but as discussed above does not give a complete reflection of the microbial community. Thus, despite the vastly superior ability of molecular techniques to identify bacterial phylotypes, species-specific identification of bacteria remains superior in

culture-based techniques [39]. For this reason, improved specificity of molecular detection techniques is necessary for not only a more complete understanding of the human sinonasal microbiome and other microbial communities, but also to be able to use this approach for decision making in the clinical context. Lastly, identifying the microbial taxa at play in different diseases with higher specificity will enable more directed experimental follow-up studies.

To take advantage of newer PacBio sequencing chemistry, improve upon data processing methods, and apply FL16S gene sequencing to a clinically relevant context, we describe a new pipeline (MCSMRT, “Microbiome Classification by Single Molecule Real-time Sequencing”). We show using two mock communities (one with 280 bacterial species) that FL16S CCS reads offer unprecedented accuracy and precision. We then explore bacterial diversity in the human nose and paranasal sinuses using results from MCSMRT, investigating not only bacterial diversity among subjects but also diversity within subjects at distinct sub-anatomical sites.

Results

Microbial community profiling by FL16S deep sequencing and CCS error correction

The taxonomic and phylogenetic resolution of microbial community profiling via 16S rRNA gene sequencing was increased by using PacBio RSII to generate FL16S sequences from mock and human sinonasal microbial communities. We combined a circular sequencing template approach with the long DNA polymerase read-lengths provided by the PacBio sequencing technology. This provided for multiple sequencing passes of each molecule, enabling the generation of CCS reads of exceptionally high quality [19, 21]. To analyze these data, we developed a new bioinformatics pipeline, MCSMRT, building upon the UPARSE pipeline [40], which (a) processes and filters PacBio CCS reads generated from multiplexed samples, (b) de novo clusters high-quality FL16S sequences into “operational taxonomic units” (OTUs), (c) taxonomically classifies each read and assigns confidence values at each taxonomic level, and (d) quantifies the abundance of each OTU based on the full CCS read dataset (Fig. 1). This processed data is suitable for downstream microbiome analyses using standard tools [41–44]. We further apply our classifier to all filtered reads and also allow for detection of amplicon sequence variants (ASVs) among groups of related sequences via minimum entropy decomposition (MED). Details are in the “Methods” section and Additional file 1, and the MCSMRT software documentation is freely available (<https://github.com/jpearl01/mcsmrt>).

Below, we demonstrate the robustness and high taxonomic and phylogenetic resolution of our experimental

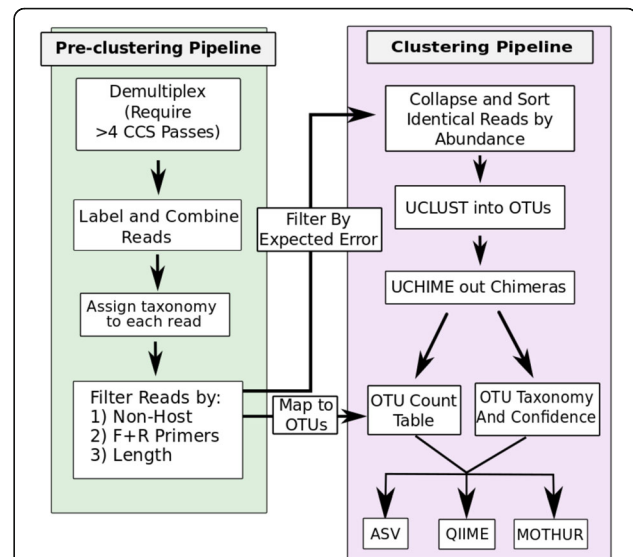


Fig. 1 Overview of the MCSMRT pipeline represented as a flowchart. MCSMRT analysis of 16S rRNA reads from the PacBio is carried out in two steps: In the pre-clustering step, CCS reads are generated during demultiplexing, labeled by sample, pooled together, and then filtered based on several criteria (length distribution, terminal matches to the primer sequences, and not aligning to a provided host or background genome sequence). Before the clustering step, CCS reads are filtered based on cumulative expected error ($EE < 1$). The clustering pipeline uses UCLUST to identify and sort unique sequences based on their abundance, clusters CCS reads into OTUs (filtering out chimeric reads during clustering), and then using uchime after clustering as a second chimera removal step. An OTU count table is created by mapping the filtered results from the end of the pre-clustering pipeline, and each OTU is taxonomically classified based on a representative “centroid” sequence. Taxonomic classification is also applied to all filtered reads, and ASV detection by MED can be applied on multiple alignments of sets of related sequencing, grouped by either OTU or binned by taxonomic level

and bioinformatics approach. We use results from two distinct mock microbial communities: one from the Biodefense and Emerging Infections Research Resource (BEI) and the other from the “Critical Assessment of Metagenome Interpretation” (CAMI) project [45]. We then applied FL16S gene sequencing to ask how the healthy human sinonasal microbiome varies among individuals and among sub-anatomical sites within individuals (Table 1; expected mock community compositions in Additional file 2: Tables S1 and S2).

CCS and filtering for FL16S reads

From each sample across the three types of communities, PCR was used to amplify FL16S genes (~1.5 kilobases [kb]) from total DNA purifications using primers that targeted conserved regions at both ends of the gene. These primers also contained terminal asymmetric barcodes to allow for pooling and subsequent demultiplexing of multiple samples into the same SMRTcell (Additional file 2:

Table 1 Community characteristics

	BEI-EC	CAMI	HSNM-MS
Source	Biodefense and Emerging Infections Research Resource Repository (BEI)	Joint Genome Institute (JGI)	Philadelphia Veterans Affairs Medical Center
Type	Mock community	Mock community	Wild community
Details	BEI Even B (Catalog ID: HM-782D)	CAMI competition	Human Sinonasal Microbiome
Number of bacterial species	20	282 (308 strains)	Unknown
Number of other species	2 (archaeal and fungal)	2 (archaeal)	Unknown
Species distribution	Reportedly even	Widely varying	Unknown
Other	Pooling DNA based on 16S qPCR	Pooling DNA based on genomic DNA mass	12 subjects, 6 sites, swab, and biopsy

Source and composition characteristics of BEI-EC and CAMI (mock communities) and HSNM-MS (Human sinonasal community)

Table S3). For the simpler BEI mock community, we tested PCR parameters by varying the polymerase (GoTaq vs. AccuPrime), the cycle number (22 vs. 35), and the presence of excess off-target DNA (i.e., 10-fold excess of genomic DNA from U937 lymphoblast lung cell line).

We used the PacBio RSII (P6-C4 chemistry) to collect 3,446,849 polymerase reads in total across the three communities (Table 2). We typically obtained ~50–60K polymerase reads per RSII SMRTcell, in which 3–4 barcoded FL16S libraries were pooled together and subsequently demultiplexed. More than half of polymerase reads were typically > 20 kb, such that the average polymerase read included ~12 complete sequencing passes around each molecule (average 1422 base pair [bp] inserts). Those polymerase reads with > 4 passes were used to generate error-corrected CCS reads, whose quality was dramatically improved (mean cumulative expected errors, EE, of 4.9 per kb) compared to polymerase reads to (EE/kb = 139.4). CCS reads with 1–4 passes had considerably lower quality (mean EE/kb = 193; 98.9% of these reads had EE > 10), and these were not considered further. The requirement for at least 5 passes resulted in a large reduction in overall yield compared to total polymerase reads but massively increased confidence in base calling (Table 2).

A series of additional filters were applied to these CCS reads to eliminate off-target sequences: (a) a size filter, (b) a filter against background (host) sequences, and (c) a primer matching filter. Collectively, these filters eliminated ~20–40% of CCS reads (Fig. 1a, Table 2, Additional file 3: Figure S1, Additional file 2: Table S2).

Size filter

CCS reads were removed if they were outside the thresholds of FL16S sequences (those between 0.5–2 kb were retained). This 2–3% of all CCS reads were mostly dimeric 16S sequences ~3 kb long, most likely created during the ligation step of library preparation (Additional file 3: Figure S2).

Host filter

CCS reads were removed if they aligned to a background genome (in this case, the human GRCh37, or hg19, reference). Notably, only 19 of ~160K reads from the BEI mock community samples mapped to the human genome, despite half of these samples including a

Table 2 Bulk sequencing and filtering stats for all three communities

Source ¹	BEI-EC ²	CAMI ³	HSMC-MS ⁴
Type	Mock	Mock	Human
Total samples	8	1	122
Total # SMRTcells	4	1	85
Total pol reads	396,625	53,164	2,997,060
N ₅₀ pol read length in kb	23,523	21,220	18,837
Avg read length in pol reads	12,997	11,232	11,940
Avg Phred-quality in pol reads	9.89	9.5	9.35
Avg EE of pol reads ⁵	1868.82	1727.44	1653.9
EE per kb of pol reads	143.79	153.8	138.52
Avg CCS passes	14.04	13.06	11.48
Avg Phred-quality in CCS reads	40.46	40.84	39.26
Avg CCS length	1454	1481	1417
Avg EE for all CCS reads ⁵	5.36	3.15	7.16
EE per kb of CCS reads	3.69	2.13	5.05
Total CCS yield (> 4 passes)	163,689	19,576	787,302
Size filtered (0.5–2 kb)	131,413	16,061	498,007
Host filtered	163,670	19,574	704,935
Primer matched	131,856	16,156	498,820
Percent passed primary filters	80.4	81.8	63.1

¹In addition to mixed species communities, seven independent negative reagent controls and four positive controls (2 x DNA from pure cultures of *Escherichia coli* and *Agrobacterium tumefaciens*) were run

²Three conditions: polymerase (GoTaq vs. AccuPrime) × PCR cycles (35 vs. 22) × excess DNA (no human DNA vs. 10-fold excess)

³Four independent libraries

⁴Twelve subjects with healthy sinuses sampled at six sinonasal anatomic locations, both swab and biopsy. Twenty samples (mostly sites E and F) were not collected or not run

⁵EE is cumulative expected error across the full read

10-fold excess human DNA (extracted from U937 human lymphoblast lung cell line to minimize contamination from the human microbiome). This indicates no appreciable off-target priming or contaminating fragments from U937 cell line DNA added at a 10:1 excess (by mass). However, samples with added human DNA had marginally lower sequencing yields (~25%), possibly indicating a weak inhibitory effect by excess off-target DNA (Additional file 3: Figure S3A, effect on CCS yield, Tukey's HSD $p = 0.048$ for 10 excess human DNA, but $p > 0.2$ for polymerase or cycle number).

By contrast, a much larger proportion of CCS reads from the human sinonasal communities mapped to the human genome (9.9%), suggesting off-target amplification of human DNA when in vast excess over bacterial DNA (alternatively, the U937 cell line DNA used for the BEI experiment may have lacked some or all off-target priming sites present in the human reference). Supporting this interpretation, biopsy samples had substantially higher total DNA yields after extraction than swabs (though no obvious differences in PCR yield, Additional file 3: Figure S4), and reads derived from human were significantly more abundant in biopsy samples (Additional file 3: Figure S3B and C, comparing biopsies and swabs, Tukey's HSD $p < 0.01$ for total CCS yield or % human contaminants, but $p > 0.8$ when for patient or site). The 105,801 reads in the sinonasal dataset that mapped to the human genome aligned to 9716 distinct genomic positions, but they were highly enriched at only a few (67.9% of mapped positions had only a single mapped read, but 58.2% of reads started at only 16 positions and had >1000-fold coverage, Additional file 3: Figure S5). These data suggest an off-target priming effect at high excess concentrations of human DNA with "hotspots" for off-target priming, along with a proportion of library molecules carrying apparently random human genomic DNA fragments. To confirm that reads mapped to the human genome were not improperly aligning true bacterial 16S genes, utax classification of all human-mapping reads showed that all had extremely low confidence assignments to the bacterial domain (<0.1), indicating a probable host origin.

Primer matching

CCS reads were required to have the forward and reverse primer sequences each found once and oriented correctly at the ends of the sequence, and this removed 12–18% of reads (Table 2, Additional file 2: Tables S4 and S16, Additional file 3: Figure S6, Additional file 1). Primer matching also served to determine the orientation of the 16S gene in each CCS read, so reads were reverse complemented when the reverse primer came first. Finally, primers were trimmed from reads. In principle, this loses several taxonomically informative sites, since the

primers contained four degenerate bases; however, in practice, the primer sequence seen in a given read was random with respect to the taxonomic source of that 16S gene. This is most easily illustrated from control sequencing of 16S rRNA genes amplified from clonal cultures of *Escherichia coli* K12 MG1655 and *Agrobacterium tumefaciens* NTL1 (Additional file 3: Figure S7 and S8).

Clustering CCS reads into OTU

Profiling the bacterial composition of a microbiome often begins by clustering sequences with high sequence identity into OTU, with a standard cutoff of 97% [46, 47], though sometimes, other cutoffs are used [48–50]. Newer approaches to grouping together related sequences avoid using similarity thresholds but instead define amplicon sequence variants (ASVs) based on controlling for variant sites arising due to sequencing error; these methods include oligotyping, minimum entropy decomposition (MED), and DADA2 [43, 44, 51]. Here, we initially show results with OTU clustering and then show how MED can further discriminate species whose 16S rRNA genes diverge by less than the threshold used for OTU picking.

To first cluster reads and identify representative "centroid" OTU sequences, we used the UCLUST algorithm [52] (Additional file 1), which filters chimeric 16S sequences by identifying apparent hybrids between distinct OTU as they accumulate in the dataset (here called CHIM1). A second chimera filter (CHIM2, using uchime) [53] then removes centroid OTU that appear to be hybrids of distinct 16S sequences in the curated Ribosomal Database Project (RDP) Gold database [54] (Additional file 1). The abundance of each OTU in each sample was then determined by counting the number of filtered CCS reads that aligned to each centroid (Fig. 1b). Though OTU clustering can collapse or separate distinctly named species into the same or different OTU, it systematically defines taxa in a uniform way that does not depend on taxonomic nomenclature [55, 56].

Sequencing error increases observed OTU counts

Erroneous base calls in CCS reads risk artificially inflating the number of OTUs, since reads with sequencing errors in similar 16S rRNA genes can be spuriously separated into distinct OTUs, especially if their actual divergence is near the 3% divergence cutoff and/or they are short. Thus, a final pre-clustering filter was applied based on cumulative expected error (EE, or the sum of error probabilities across all positions in a read as determined from Phred-scaled base quality scores). This measure has previously been shown to discriminate against error-prone sequences better than the average quality score [57].

Analysis of the BEI mock community

Using the BEI mock community to examine the relationship between CCS passes and EE, we found, as expected, that reads with more CCS passes had a lower median EE (Fig. 2, linear model of $\log(\text{EE})$ vs. CCS passes gives $R^2 = 0.22$). 98.9% of reads with less than five CCS passes had $\text{EE} > 10$ (across ~ 1.4 kb total length) and were not considered further (Fig. 2).

To empirically determine an appropriate EE cutoff for clustering CCS reads into OTUs, we compared the expected number of OTUs in the BEI mock community to that obtained by OTU clustering at different EE cutoffs. The expected number of OTUs was 19, since two of the 20 species' 16S rRNA genes differed by at most 23 nucleotides (*Staphylococcus aureus* and *Staphylococcus epidermidis* have only 1.4% divergence, less than the OTU clustering cutoff of 3% divergence, Additional file 2: Table S1). As expected, increasing the stringency of the EE filter reduced the total number of CCS reads available for OTU clustering, as well as the total number of OTUs detected (Fig. 3a). Using $\text{EE} \leq 1$ (one or fewer expected errors per read) retained less than half (40.1%) of filtered CCS reads, but these clustered into the 19 OTUs expected. Decreased stringency (higher EE cutoffs) increased the total OTUs detected, dramatically for cutoffs of $\text{EE} \leq 8$ and above; using no expected error threshold ($\text{EE} \leq 128$), 3453 OTUs were detected, more than 100-fold greater than the true number. In summary, using a high stringency expected error cutoff for OTU clustering reduced the number of reads available

for clustering but provided exact total OTU counts for the BEI mock community.

Sequence length and OTU clustering

To examine how OTU clustering would be affected by using partial instead of FL16S gene sequences for the BEI mock community, we performed in silico primer matching and trimming on the full-length CCS reads for three short-read primer pairs commonly used for microbial community profiling, namely primers targeting the V1-V3, V3-V5, or V4 hypervariable regions of the 16S rRNA gene, which are amenable to analysis using Illumina short-read sequencing (Additional file 2: Table S3). Using these in silico short-read data produced dramatically higher OTU counts than predicted, even when using substantially more stringent EE cutoffs (Fig. 3b). Thus, for example, whereas full-length V1-V9 reads clustered into the expected 19 OTUs at $\text{EE} \leq 1$, reads truncated to include only the V3-V5 hypervariable regions (average = 536 nt) clustered into 116 OTUs. The other truncated sequences (V1-V3 and V4) had even higher elevated total OTU counts. Even when the shorter length of partial sequences was compensated for by using an eightfold higher stringency cutoff ($\text{EE} \leq 0.125$), spurious OTUs were still detected, e.g., 58 OTUs were detected with the V3-V5 truncated reads, substantially higher than expected. The comparisons above relied on in silico truncation of full-length reads from the same PacBio dataset to maintain consistent error profiles, but inflated OTU counts have also been reported in published results for truncated 16S from

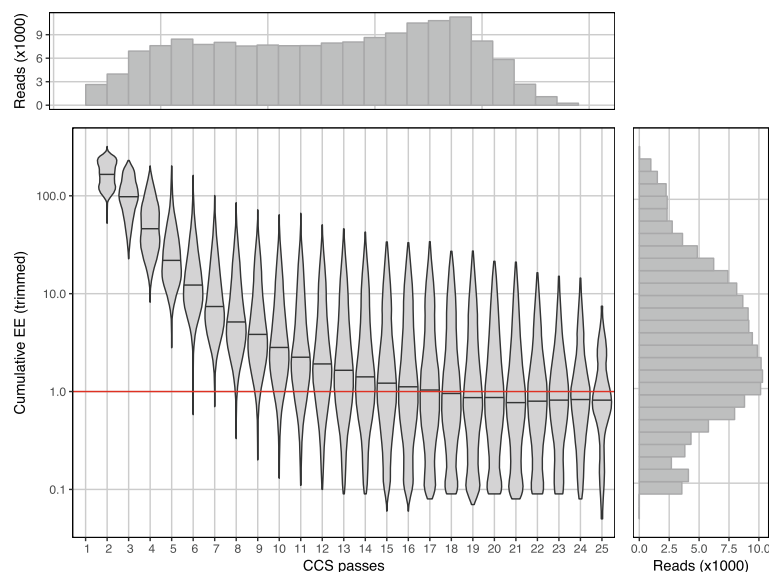


Fig. 2 Distribution of reads at different CCS passes and cumulative expected error values (EE) in the BEI mock community. Violin plot showing the distribution of cumulative EE (after primer matching and trimming) at different CCS passes. Reads with less than two CCS passes were not reported by PacBio CCS software. Histograms at the top and right show read count by CCS and EE, respectively. The 35 reads with 26 to 46 CCS passes are not shown (median EE = 0.22). Subsequent analyses used only CCS reads with > 4 passes

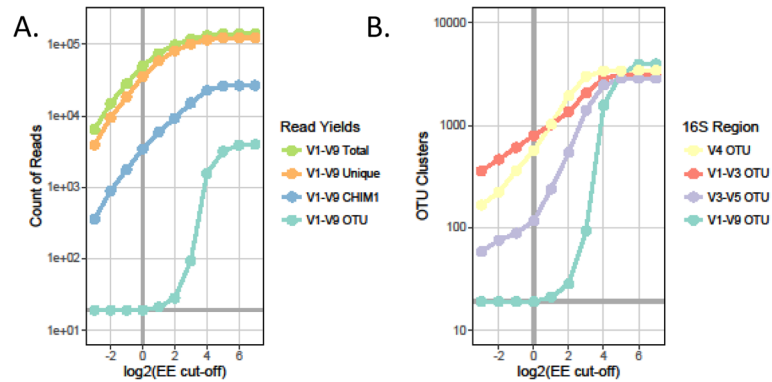


Fig. 3 Clustering of post-filtered CCS reads into OTUs. **a** Count of total, unique, CHIM1, and centroid OTU reads at different maximum EE thresholds. **b** Count of total OTU detected using full-length or truncated reads at different maximum EE thresholds

the BEI mock community collected using both 454 pyrosequencing and Illumina MiSeq short-read technology [58–62]. For comparison, we also applied closed-reference OTU clustering at a 97% cutoff via QIIME2 and found inflated OTU counts using either FL16S (115 OTUs) or V3-V5 truncated reads (202 OTUs) (if counting only OTUs with > 4 reads mapping, then FL16S detected 60 OTUs and V3-V5 detected 115 OTUs), though FL16S still detected fewer OTU counts than truncated reads. Finally, we similarly found elevated OTU counts in a re-analysis of Illumina MiSeq data for the V1-V3 region of the same BEI mock community [61] through our pipeline, finding 171 OTUs at $EE < 1$ (40 OTUs at $EE < 0.125$). This suggests that inflated OTU counts when using partial 16S sequences are independent of the specific PCR conditions or the particular error profile of PacBio CCS reads.

The above results underline the value of using FL16S to minimize the effect of sequencing errors on de novo OTU cluster counts. They also indicate that methods that profile taxonomic composition using partial 16S rRNA genes may be prone to overestimating bacterial diversity. For all subsequent analyses, we used only CCS reads with $EE \leq 1$ for OTU clustering, and then mapped all reads passing all pre-clustering filters onto these centroid OTU's to obtain abundance data (Pre-clustering Pipeline, Clustering Pipeline, Additional file 1).

Taxonomic classification of FL16S reads

Bacterial taxonomic nomenclature has traditionally been based on physiological and other microbiological traits (e.g., virulence) rather than 16S rRNA gene sequences, so the accuracy and precision with which a read can be taxonomically defined is dictated by a combination of organism-specific criteria used for naming species, the quality and completeness of the database used, and the distribution of informative variable sites within the 16S rRNA among named taxa [9]. Unfortunately, commonly used databases for classifying 16S rRNA gene sequences,

namely RDP and Silva [54, 63, 64] (although see [65] for a new way of extracting FL16S reads with species-level classifications from RDP), do not provide species-level taxonomic identifiers [66]. Another problem with these databases is the absence of representative sequences from genera present in our mock communities (for example, the genus *Clostridium* was not found in RDP). Though 285,289 sequences in the popular Greengenes database do have species labels, only 631 of these are unique species. Although Greengenes ($N = 1,262,986$) and Silva ($N = 1,922,213$) have vastly more taxonomically classified sequences than RDP ($N = 8978$), in part because they computationally assign taxonomies to sequences from environmental microbiome surveys [63, 67], most of the sequences in these databases are of partial length 16S rRNA genes. While these databases are appropriate in many cases, we needed to make a database of FL16S sequences with species-level taxonomic information.

To classify CCS reads (including centroid OTUs and MED node representatives) based on bacterial taxonomy to the species level and also provide confidence values at each taxonomic level, we trained a utax classifier on a custom-built database of FL16S gene sequences downloaded from NCBI (16S rRNA Microbial Database, Additional file 1). Most FL16S sequences available at NCBI ($N = 17,764$) could be associated with a taxonomic ID (txid) using the gid accession number, allowing us to extract, parse, and configure sequences in the database to create a utax-compatible Linnaean hierarchy that included 11,055 distinctly named species spanning 367 bacterial families (Additional file 1). The number of distinct families present in this NCBI database was 367, whereas RDP had 366, Silva had 302, and Greengenes had 514. We recognize that other researchers may not prioritize species-level taxonomic assignments and instead favor high breadth. To that end—because MCSMRT is based on the UPARSE pipeline—any correctly formatted database may be used in place of our custom

one. UPARSE-formatted databases for Greengenes, Silva, and RDP are available and may be found at [68].

We next generated and compared the accuracy of utax classifiers built from full-length or partial V3-V5 16S rRNA genes by classifying the database sequences themselves. In this context, incorrect classification could arise in particular due to distinct named species with highly similar sequences. The full-length classifier gave an incorrect label only 1.0% of the time ($N = 173$ mistakes), compared to 13.2% of the time using the truncated classifier ($N = 2295$ mistakes). Indeed, when the two classifiers disagreed, the full-length call was much more frequently correct (2×2 contingency table: 15,081 both correct, 2137 only full-length correct, 15 only truncated correct, and 158 neither correct). Furthermore, species-level confidence values were higher 81.3% of the time using the full-length classifier (mean 81.7%, median 92.7%) compared to the truncated classifier (mean 71.1%, median 82.7%). These results show the value of using full-length compared to partial 16S gene sequences for accurate taxonomic assignment.

The assignment made for each CCS read is associated with confidence values at each taxonomic level, and low values could arise for several reasons aside from the quality of the sequence data. In particular, sequences labeled as a distinctly named species could have other equally good matches, or nearly so. In order to determine what species might end up assigned to a particular centroid OTU read, we clustered the NCBI database sequences (17,776 in total, 99.1% unique) at the same threshold level (97% identity), thereby grouping species belonging to the same “database OTU” (dbOTU). Since UCLUST relies on abundant unique sequences to initiate centroids and also drops putative chimeric sequences during clustering, we instead applied hierarchical clustering (average linkage, using pairwise percent identity values from all-by-all blast, and separating dbOTU clusters at a 3% difference level). This method is unaffected by the order of the sequences and included all database entries.

Hierarchical clustering of NCBI sequences resulted in 6065 dbOTU, of which 66.9% of clusters had a single species (93.2% had a single genus), whereas 14.6% of clusters had the same species split over more than one dbOTU (Additional file 4). Some dbOTUs consisted of many species. For example, the top three most species-rich dbOTUs collectively contained 453 distinctly named *Streptomyces* species indicating that 16S rRNA clustering at 3% divergence poorly discriminates among named species in this genus [69]. These results reflect the variability with which different bacterial taxa are named compared to how they group based on divergence in their 16S rRNA [49, 50] (Additional file 3: Figure S9). Collectively, clustering the FL16S gene sequences from NCBI indicated that

assignments of individual CCS reads, OTU centroids, or MED node representatives will result in high-confidence species-level classification for high-quality FL16S gene sequences; imprecision due to distinct species belonging to the same dbOTU will affect about a third of 16S rRNA sequences in the database, but these can be flagged by low confidence values from the utax classifier and by cross-referencing to dbOTU clusters to identify other possible “nearly best hits.”

Nearly all reads collected from reagent controls failed filtering and classification steps

In advance of the studies described above, we pre-screened multiple PCR reagents and DNA polymerases to identify those that produced no observable amplification when using reagent controls. In addition, seven negative control reagent samples were sequenced in parallel with mock community and sinonasal samples. After demultiplexing, these controls produced a total of only 54 reads. Of these, only three reads passed the filtering criteria described above. Taxonomic classification on all 54 reads returned only five reads with >10% bacterial domain-level confidence values (all five gave 100% confidence). These five reads had genus-level classification as *Finegoldia* (95.3% confidence), *Propionibacterium* (16.6%), *Streptococcus* (22.7%), and two as *Staphylococcus* (0.8%). Four of these reads had 0% species-level confidence, while *Finegoldia magna* had 94% species-level confidence. These results demonstrate that our laboratory and bioinformatics methods produce extremely low levels of contamination from off-target bacterial nucleic acids generated from our reagents. Further consideration of controlling for contaminants would require direct empirical measurements of 16S copy number in each sample [70].

BEI mock community composition

OTU classification

The 19 distinguishable OTUs in the BEI mock community were readily identified and accurately classified (Table 3). All but three of the centroid OTU—including three distinct *Streptococcus* species—were correctly classified to the species level. The three discrepancies were, however, reflected by low confidence values assigned by utax, as well as by clustering of discrepant taxa into the same dbOTU.

All 19 centroid OTUs were correctly classified to the family level, with one genus-level discrepancy: classification of *Escherichia coli* as *Shigella flexneri*. This incorrect assignment is not a surprise; indeed, the matching dbOTU contained 54 sequences that were assigned to *Escherichia*, *Shigella*, *Citrobacter*, and *Salmonella*, all genera known to have low levels of divergence among their 16S rRNA genes [71, 72]. Two species-level assignments were incorrect: (a) *Bacillus cereus* was classified

Table 3 Taxonomic classification of BEI

Centroids full-length	Centroids v3-v5			Clustered from V3-V5 truncated reads (Top 19)			Mothur/Silva full-length (Top 19)				
	Genus conf	Species conf	Genus conf	Species conf	Genus conf	Species conf	Genus conf	Species conf	Genus conf		
<i>Acinetobacter_baumannii</i>	0.98	0.94	0.88	0.42	0.8784	0.42	Acinetobacter_baumannii	0.8784	0.42	Acinetobacter	100
<i>Actinomyces_odontolyticus</i>	0.99	0.94	0.97	0.62	0.9728	0.62	Actinomyces_odontolyticus	0.9728	0.62	Actinomyces	100
<i>Bacillus_anthraxis</i>	0.99	0.49	0.96	0.42	0.963	0.31	Bacillus_cereus	0.963	0.31	Bacillus	100
<i>Bacteroides_vulgatus</i>	0.99	0.94	0.98	0.90	0.9847	0.90	Bacteroides_vulgatus	0.9847	0.90	Bacteroides	100
<i>Clostridium_beijerinckii</i>	0.99	0.72	0.97	0.42	0.965	0.31	Clostridium_roseum	0.965	0.31	Clostridium_sensu_stricto	100
<i>Deinococcus_radiodurans</i>	0.99	0.94	1.00	0.90	0.9994	0.90	Deinococcus_radiodurans	0.9994	0.90	Deinococcus	100
<i>Enterococcus_faecalis</i>	0.99	0.94	0.88	0.83	0.8784	0.83	Enterococcus_faecalis	0.8784	0.83	Enterococcus	100
<i>Shigella_flexneri</i>	0.69	0.49	0.67	0.42	0.6072	0.31	Shigella_flexneri	0.6072	0.31	Escherichia_Shigella	100
<i>Helicobacter_pylori</i>	0.99	0.93	0.98	0.90	0.9787	0.90	Helicobacter_pylori	0.9787	0.90	Helicobacter	100
<i>Lactobacillus_gasseri</i>	0.99	0.72	0.97	0.42	0.9728	0.42	Lactobacillus_gasseri	0.9728	0.42	Lactobacillus	100
<i>Listeria_monocytogenes</i>	0.99	0.72	0.96	0.42	0.961	0.31	Listeria_ivanovii	0.961	0.31	Listeria	100
<i>Neisseria_meningitidis</i>	0.97	0.93	0.95	0.62	0.9533	0.62	Neisseria_meningitidis	0.9533	0.62	Neisseria	100
<i>Propionibacterium_acnes</i>	0.97	0.96	0.97	0.91	0.965	0.91	Propionibacterium_acnes	0.965	0.91	Propionibacterium	100
<i>Pseudomonas_aeruginosa</i>	0.97	0.93	0.93	0.62	0.9329	0.62	Pseudomonas_aeruginosa	0.9329	0.62	Pseudomonas	100
<i>Rhodobacter_sphaeroides</i>	0.97	0.60	0.92	0.42	0.9228	0.31	Rhodobacter_sphaeroides	0.9228	0.31	Rhodobacter	100
<i>Staphylococcus_epidermidis</i>	0.99	0.72	0.67	0.42	0.6661	0.31	Staphylococcus_epidermidis	0.6661	0.31	Staphylococcus	100
<i>Streptococcus_agalactiae</i>	0.99	0.95	0.97	0.90	0.9728	0.90	Streptococcus_agalactiae	0.9728	0.90	Streptococcus	100
<i>Streptococcus_mutans</i>	0.99	0.95	0.98	0.90	0.9768	0.90	Streptococcus_mutans	0.9768	0.90	Streptococcus	100
<i>Streptococcus_pneumoniae</i>	0.99	0.49	0.98	0.42	0.9768	0.42	Streptococcus_pneumoniae	0.9768	0.42	Streptococcus	100

Classification results from the BEI mock community. The performance of MCSMRT FL16S centroid OTU assignments (column 1), centroids truncated to V3-V5 (column 2), top 19 centroids found when clustering V3-V5 truncation reads (column 3), and the top 19 OTUs identified using Mothur/Silva with FL16S reads (column 4)

as *Bacillus anthracis*; these two differ by only two nucleotides in their 16S genes and share the same dbOTU with eight other *Bacillus* species; and (b) *S. aureus* and *S. epidermidis*—whose 16S genes differ by 23 nucleotides—were collapsed into a single OTU, with the centroid called *S. epidermidis* (the matched dbOTU consisted of 39 additional staphylococci, though see below). Species-level confidence values were 0.49 and 0.71 respectively. By contrast, only one correctly classified OTU, *Streptococcus pneumoniae*, had a species-level confidence value < 0.50 and this species shared the same dbOTU with 13 additional *Streptococcus* species.

Truncation of CCS reads prior to clustering considerably worsened classification; in addition to increasing the number of total OTU, CCS reads truncated to their V3-V5 region prior to clustering resulted in seven misclassified OTU among the top 20 most abundant OTU (Table 3). To isolate the effects of truncation on classification alone, rather than both clustering and classification, we also truncated the 19 centroid OTU to their V3-V5 region and classified these using a utax classifier built from a database of sequences also truncated to V3-V5. This showed reduced species-level confidence values but also more miscalled taxa (Table 3). These results show that use of FL16S gene sequences provides substantially improved taxonomic identification of centroid OTU compared to truncated 16S rRNA sequences, to the extent that bacterial nomenclature allows.

For comparison, we used the Greengenes v13_8 database to classify closed-reference OTU identified by QIIME2. The resulting classifications were typically at higher taxonomic levels and more often incorrect. For example, using V3-V5 sequence, five of the top 19 most abundant closed-reference OTUs had species level annotation (*Streptococcus agalactiae*, *Listeria seeligeri*, *Clostridium paraputrificum*, *Staphylococcus saprophyticus*), but only one of these matched a species in the BEI community (*S. agalactiae*). Greengenes classification with FL16S made only two classifications to the species level, both incorrect (*S. saprophyticus* and *Alkanindiges illinoisensis*). Often classifications were to much higher levels (e.g., family Planococcaceae or even the domain bacteria). Taxonomy results using QIIME2 are reported in Additional file 2: Tables S5–S8.

We also applied the Mothur pipeline to our primer-matched EE-filtered CCS reads from the BEI mock community, following the recommended settings reported in Schloss et al. [23]. As previously reported, Mothur called numerous spurious low-abundance OTUs ($n = 1197$), but the top 19 most abundant OTUs were correctly classified to the expected 19 genera in the BEI mock community (Table 3 and Additional file 2: Table S9). We note that only genus-level assignments were reported, and all of these were assigned a 100% confidence level, suggesting that this algorithm may overestimate confidence when using FL16S.

Finally, we compared the performance of MCSMRT to FL16S CCS reads that were previously collected from a distinct mock community and processed using Mothur reported in Singer et al. [24]. MCSMRT with default settings identified 22 OTUs (of 23 expected, three additional OTUs were detected but were assigned the same species-level taxonomy as another OTU) and correctly classified all 22 of them to an expected member of this mock community (Additional file 2: Table S10). One low abundance taxon (*Nocardiopsis dassonvillei*) was not detected by MCSMRT, likely because the few reads collected failed our filtering steps. By contrast, Singer et al. also identified numerous spurious low-abundance OTUs, arriving at the correct number only after first removing these. Notably, although appropriate for purposes of their error analysis, the Singer et al. study assigned CCS reads to mock community member by mapping reads to reference sequences of those expected community members, whereas MCSMRT clustered reads into the correct groups and labeled them correctly without knowing what the expected community composition was.

Relative abundance and sequencing error

The abundance of each OTU was estimated by assigning all filtered CCS reads (with no EE threshold) to a centroid OTU with a maximum of 3% divergence for a hit to be counted. The 20 bacterial species in the BEI mock community were expected to have equimolar abundances of their 16S rRNA genes, and for most species, we detected a roughly even mock community composition for most species (Additional file 3: Figure S10). Several were outliers: (a) the two *Staphylococcus* species were binned together as *S. epidermidis* (as described above); (b) *Bacteroides vulgatus* and *Helicobacter pylori* were over-represented, especially at high PCR cycle number; and (c) five species were found at lower than expected abundances across PCR conditions.

Independent analyses of the same BEI mock community by Illumina MiSeq for V3-V5 have previously found the same taxa elevated or depleted, suggesting that these taxa actually are at unequal concentrations in this mock community [73, 74]. The primers we used have perfect identity with all BEI bacterial strains' reference 16S rRNA gene sequences and are distinct from the Illumina-based analyses, so the compositional biases seen are not likely to be due to primer choice or PCR conditions [61].

We next evaluated the impact of (a) chimeric sequences on relative abundance measurements and (b) “true” substitution errors. First, all CCS reads were run through UCLUST with no filters other than requiring > 4 CCS passes to identify likely CHIM1 chimeras, and then all CCS reads were aligned to the 16S reference sequences

from the BEI community to determine their likely source. This again found relatively even abundances for each taxon with the exception of those mentioned above. Increased cycle number also increased the variance among taxa in their relative abundances, but the inclusion of chimeric reads had little effect (Additional file 3: Figure S11). Second, the number of base mismatches (excluding gap characters) was calculated between each non-chimeric read and its most similar reference sequence, estimating the number of “true” substitution errors made during sequencing. Notably, intragenomic variation in 16S rRNA gene sequences [75] also contributes to putative substitution errors. This analysis indicates that the AccuPrime polymerase made fewer errors than GoTaq polymerase but that errors made by either polymerase were insufficient to inflate OTU numbers when using full-length sequence (Additional file 3: Figure S12). Overall, the mean error rate across all BEI mock community reads was (a) for $EE \leq 1$, 6.05 substitutions per read and 2.00 indels per read, and (b) for $EE > 1$, 14.05 substitutions per read and 6.86 gaps per read. However, examination of these error rate distributions for several taxa (Additional file 3: Figure S12A) suggested that the single reference FL16S sequence we had for the BEI mock community may not have been identical to the strains used to build the community (e.g., for *E. coli* matched reads, almost no reads were an exact match) or that intragenomic variation had a large impact on these estimates.

Sequencing error in *E. coli* positive control CCS reads

To more precisely estimate sequencing error by ensuring we had correct reference sequences and accounting for intragenomic variation among 16S rRNA genes, we investigated sequence variation in FL16S CCS reads collected from our lab stock of *E. coli* K12 MG1655 as a monoculture positive control sample. We first obtained a finished circular assembly of our lab’s strain by shotgun sequencing on the PacBio RSII, and we identified and extracted seven FL16S genes (two were identical, but the others all differed slightly). This allowed us to obtain more confident estimates of the true error rate in individual CCS reads by globally aligning all 8038 primer-matched CCS reads to their closest matching FL16S copy from the reference genome. The 1445 CCS reads with $EE \leq 1$ had considerably lower “true” error rates (mean mismatches = 3.0 per read, mean gaps = 0.8 per read), compared to those with $EE > 1$ ($n = 6593$, mean mismatches = 5.5 per read, mean gaps = 6.0 per read). This illustrates an especially dramatic loss of errors due to indels after EE filtering, though mean errors overall still exceed the expected value of ≤ 1 , suggesting that quality scores after CCS processing are somewhat inflated. Detailed error rate statistics are found in Additional file 2: Table S11, and histograms of substitution

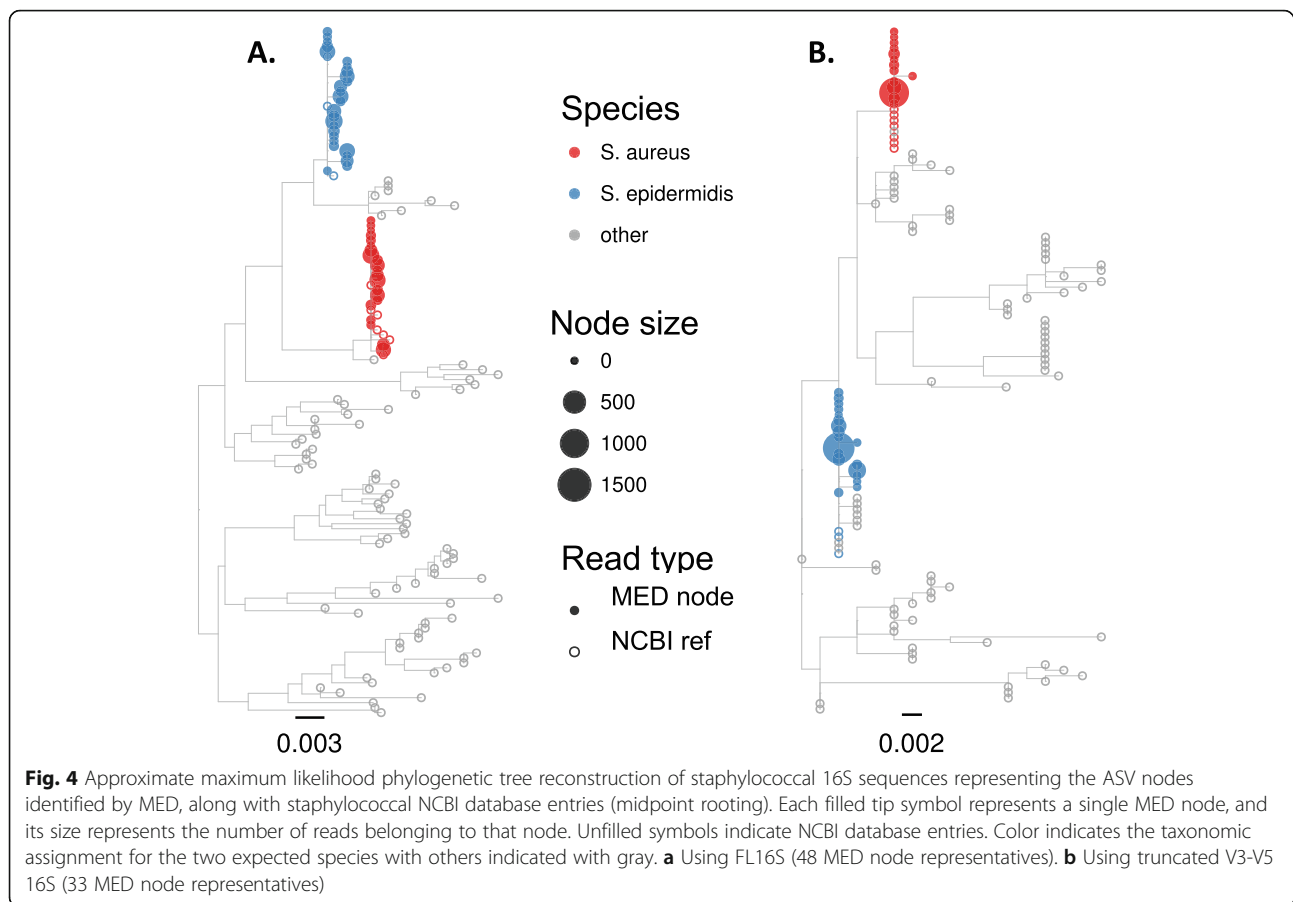
and indel errors per read are in Additional file 3: Figure S13. Additionally, we used MED to identify ASVs from multiple alignments of *E. coli* positive controls; although MED (and DADA2) can use unaligned data, they require that indels be biologically meaningful. The resulting MED node representatives were aligned with the 16S genes identified from our whole genome assembly, and an approximate ML tree shows that ASVs correctly segregated with individual genomic 16S copies (Additional file 3: Figure S13B).

Discriminating among closely related sequences

Because using FL16S should increase the number of taxonomically and phylogenetically informative sites, we reasoned that species whose 16S genes differ by less than the OTU clustering threshold would be more easily separated with full-length versus truncated 16S gene sequences. Although the two clinically important *Staphylococcus* species in the BEI mock community belonged to the same dbOTU (along with 40 other staphylococcal species, and two additional genera) and were not separated during de novo OTU clustering, they were readily distinguishable in several ways.

First, direct classification of primer-matched CCS reads from the BEI mock community identified only *S. aureus* and *S. epidermidis* among those classified to the staphylococci (1649 from *S. aureus* and 2501 from *S. epidermidis*; using unfiltered CCS reads yielded 0.71% classified to five additional staphylococcal species in 32 reads). Thus, direct taxonomic classification correctly identified both species in roughly equal proportions.

Second, we applied MED to identify ASVs for all primer-match CCS reads with $EE \leq 1$ that had been assigned to the *Staphylococcus* OTU, and the MED node representatives were taxonomically classified (Identification of Amplicon Sequence Variants by Minimum Entropy Decomposition, Additional file 1). It was necessary to apply MED to closely related sequences, because alignments of FL16S that included diverse bacteria had enough gap characters to make alignments several times longer than any one FL16S sequence. MED reduced 3171 CCS reads to 48 nodes (ASVs), and representative sequences from each ASV were used to build phylogenetic trees, also including all NCBI entries for the staphylococci (Fig. 4). The results with FL16S sequences show that *S. epidermidis* and *S. aureus* are clearly separated from the other staphylococci, as expected, and node representatives classified as each species formed clear monophyletic groups (Fig. 4a). By contrast, building trees from ASVs identified from the V3-V5 truncated reads did not clearly distinguish among staphylococcal species, and although in this case the node representatives were correctly classified, NCBI entries of other species were intermixed with those of the two expected species (Fig. 4b).



These results show that distinguishing among closely related organisms—even when sequence differences are insufficient to separate these into distinct OTUs—is strongly facilitated by use of FL16S gene sequences, and especially powerful when combined with an ASV detection method.

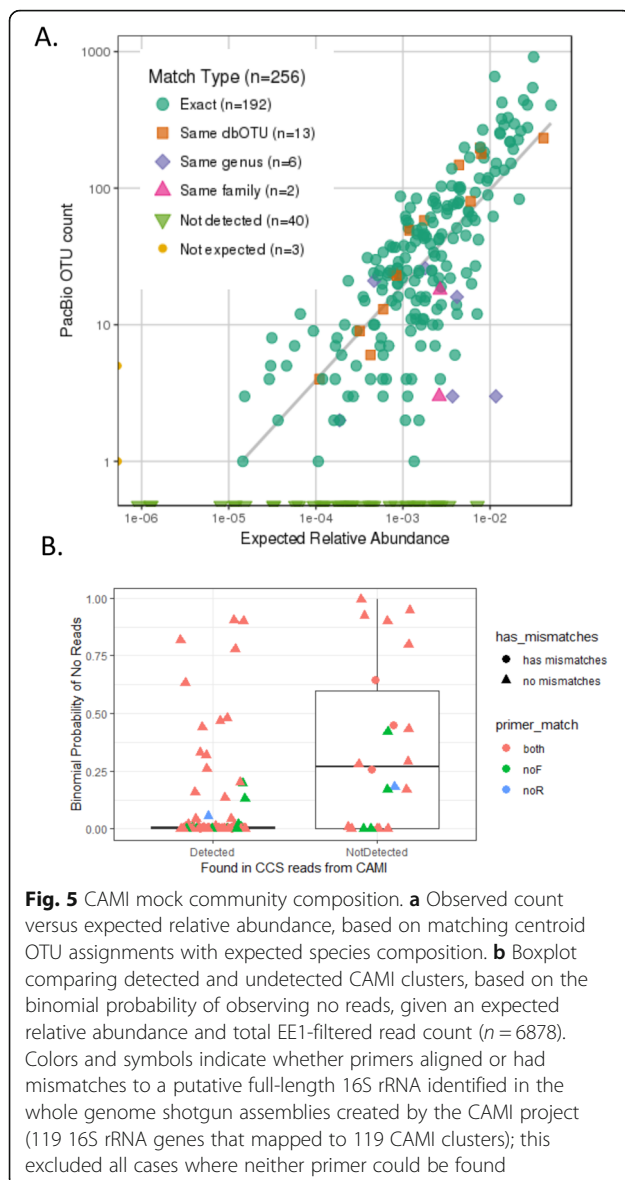
CAMI mock community composition

OTU classification

Because a curated set of FL16S gene sequences was not available for the 280 unique bacterial species present in the CAMI mock community, we first cross-referenced the expected bacterial composition (Additional file 2: Table S2) with the FL16S gene sequences in the NCBI database, finding one or more full-length sequences for all but three species, for which a taxonomy was available but not a corresponding 16S rRNA gene sequence (*Mameliella alba*, *Fusobacterium naviforme*, and *Promicromonospora flava*). In addition, 11 species names used by CAMI and NCBI were synonyms, due to revisions in species names (particularly members of the *Clostridiales* family, Additional file 2: Table S12). By cross-referencing the expected species with the NCBI dbOTUs, we found the 280 CAMI species would cluster into 253 OTUs at

the 3% divergence level, each associated with a distinct dbOTU (except the three missing species). For example, five species from the genus *Prauserella* (*Prauserella aidingensis*, *Prauserella alba*, *Prauserella flava*, *Phalophila halophila*, and *Phalophila salsuginis*) were in the CAMI community, and although there was an instance of each of those species in the NCBI database, none of these NCBI 16S sequences differed from each other by greater than 3%. Therefore, the expectation was a single OTU associated with the genus *Prauserella*, which was indeed the result. Most of these expected clusters had three or fewer named species in their corresponding dbOTU (84.4% of CAMI clusters), but they collectively comprised 586 distinct species calls in the NCBI database. This clustering allowed us to cross-reference the centroid OTUs to members of the CAMI community and identify ambiguities in the extant taxonomic classification (Additional file 2: Table S12).

FL16S gene sequencing by PacBio had exceptionally high specificity and sensitivity for identifying the bacterial constituents within the complex CAMI mock community (Fig. 5). The 16,156 filtered CCS reads (final yield from one PacBio SMRTcell) clustered into 227 OTUs (using 6878 reads at $EE \leq 1$) with 216 unique species names. Of



these, 192 centroid assignments perfectly matched up with an expected cluster in the CAMI community, thus giving 89% exact species-level matches with the centroid OTU. Nineteen more centroids could be connected to CAMI-defined clusters, either via a dbOTU (13 centroids) or a correct genus-level assignment (six centroids).

Of the five remaining OTUs detected by PacBio that did not belong to a CAMI cluster, two could be accounted for with family-level matches to CAMI clusters (the Rhodobacteraceae *Mameliella alba*, absent from the NCBI database, was classified as *Paracoccus versutus*, and *Promicromonospora flava* was identified as *Isosporicicola variabilis*). This left only three “false positive” OTUs, which accounted for a total of seven CCS reads. Two of these misidentifications were species belonging to families represented in the CAMI

community—the Rhizobiacean *Agrobacterium larrymoorei* (1 CCS read) and the Lachnospiracean *Moryella indoligenes* (five CCS reads)—and one was not represented (the Moraxellacean *Acinetobacter septicus* had a single CCS read).

In summary, with a single SMRTcell, we detected 84.2% of taxa predicted to be present with 95% species-perfect identification, and nearly every single OTU identified by FL16S gene sequencing could be attributed to a member of the CAMI community. By contrast, truncating CCS reads (or their centroid OTUs) to the V3-V5 region was less accurate and showed lower species-level confidence values, as seen above for the database sequences themselves (above and Additional file 2: Table S13).

CAMI species undetected by MCSMRT

Forty expected CAMI clusters were not identified among the centroid OTU (i.e., “false negatives”). This was, at least in part, due to under-sampling: The relative abundance of 16S rRNA genes for most CAMI species was expected to be very low (some well below our limit of detection), and all missing CAMI members had expected abundances of <1% (Fig. 5). We estimated the probability that we would observe zero reads using the binomial distribution (of 6878 EE-filtered reads used for OTU clustering) given each CAMI cluster’s expected relative abundance. This found that 15 of the 40 missing clusters (37.5%) had a <5% chance of being missed. In contrast, for the 197 CAMI clusters that we did detect, 133 had a <5% chance of being missed (67.5%), indicating that many of our false negatives are likely to be due to undersampling.

Another potential reason for failure to detect specific taxa would be failed amplification due to variation from our universal primers from CAMI community members. Because the vast majority of sequences in our NCBI database lacked our FL16S primers due to trimming, we explored possible primer matching problems by curating a set of reference sequences, first by identifying 16S rRNA genes using rnammer from Illumina whole genome shotgun assemblies created for the CAMI project. We were able to annotate 291 16S rRNA genes from 304 assemblies, and we reduced this to a set of 197 sequences representing 197 CAMI clusters. Alignment of the forward and reverse primer sequences to these found only 101 for which both primers aligned (13 had neither primer align), but many of these genes were shorter than expected for FL16S, likely because Illumina assemblies may often include incompletely assembled rRNA genes and operons. Applying a minimum length of 1400 bp reduced the set to only 119 sequences, but both primers aligned to 100 of these (17 had only reverse primer alignments, two had only forward primer alignments, and none had neither primer match).

Of the 40 “false negative” accessions, 22 of the undetected CAMI clusters were in this set of 119 FL16S reference sequences, and this included seven of the 15 identified above as unlikely to have been missed. The expected number of OTU counts (out of 6878) was still relatively low, ranging from ~5–24 reads, but especially for the two for which we expected about 24 counts (*Halobacterium salinarum* and the *Williamsia* cluster), the binomial probability of seeing no reads was $< 10^{-10}$. The *H. salinarum* sequence was one strain that had no forward primer match. A distinct *Williamsia* species that belong to its own CAMI cluster was detected, raising the possibility that these missing reads ended up clustering with this other species. Finally, three species with intermediate chances of being missed had mismatches in their primer alignments, as seen in Fig. 5b).

In conclusion, most of the missing “false negative” taxa we observed were due to under-sampling. We found some weak evidence for potential primer issues, but these were complicated by difficulties building a comprehensive reference set that included the primer sequences themselves. We also cannot rule out that some accessions may not be represented at their expected relative abundances.

Relative abundances in the CAMI mock community

The expected relative abundance of each species in the CAMI mock community was accurately reflected by the number of reads assigned to each CAMI centroid OTU by usearch (Fig. 5, Additional file 2: Table S12). The expected relative concentration of each species' 16S rRNA genes in the CAMI community was calculated using (a) the genome size estimated by from de novo assembly of shotgun sequence collected from each CAMI accession [76] (Additional file 2: Table S12) and (b) an estimate of 16S rRNA gene copy number using rrnDB [77]. For CAMI species missing from the rrnDB, the lowest Linnaean rank with members of the database was determined, and the average 16S rRNA copy number of all species under that rank was used (Additional file 2: Table S12). Remarkably, we observed a strong linear fit between observed and expected abundances (Fig. 5, $R^2 = 0.63$), showing that we not only accurately identified the species present by centroid OTU, but also accurately quantified their relative abundances, despite the low expected relative abundance of most species' 16S rRNA genes.

Phylogenetic discrimination of species in the same genus within the CAMI community

Because the CAMI mock community included 45 multi-species genera (9 with > 3 species), we next asked whether FL16S reads discriminated among species in the same genus better than truncated V3-V5 reads. We collected all filtered CCS reads that had been classified to a given multi-species genus and produced phylogenetic

trees from multiple sequence alignments of each genus-specific read set (39 multi-species genera with at least five filtered CCS reads).

Using these genus-level trees, we next assessed whether the utax-assigned species labels for each read formed monophyletic clades using MonoPhy [78] (Additional file 2: Table S14). For most genera—where the species were sufficiently diverged—trees built from either FL16S or V3-V5 truncated reads performed comparably: For 28 of 39 genera, all assigned species labels were monophyletic using either FL16S or V3-V5 reads. Examples of well-resolved genera with either marker gene length included *Clostridium* and *Desulfovibrio* (Additional file 3: Figure S14). Five more genera were non-monophyletic for an equal number of species using either marker gene length; some of this is likely due to poorly resolved species nomenclature. Examples include the genera *Azotobacter* and *Nonlabens* (Additional file 3: Figure S15). For the remaining six multi-species genera, phylogenies built from FL16S reads showed higher monophyletic grouping of species-level classifications than trees built from V3-V5 truncated reads. Two prominent examples were closely related species within the *Algoriphagus* and *Salegentibacter* (Additional file 3: Figure S16). In *Algoriphagus*, *Algoriphagus yeojeoni* appears polyphyletic for V3-V5 only, and in *Salegentibacter*, *Salegentibacter salegens* fails to resolve from *S. salinarum*. These results further demonstrate the utility of increasing the length of marker gene sequencing to capture more informative sites, thus improving phylogenetic resolution of distinct but closely related members of microbial communities.

The composition of the human sinonasal bacterial microbiome

Rhinosinusitis affects 16% of the US population [79] and accounts for one in five antibiotic prescriptions to adults in the USA in the outpatient setting, making it the most common diagnosis for outpatient antibiotic use in the USA. [80]. Thus, a more complete understanding of the resident microbial community of the upper respiratory tract is paramount to improved therapeutic interventions and reduction of inappropriate antibiotic prescriptions. Thus, we applied our bacterial microbiome profiling method to the human sinonasal cavity. We obtained samples from 12 subjects undergoing pituitary gland adenoma removal, utilizing the sinonasal cavity as a surgical corridor for access to the gland. None of the total 12 patients examined had objective or subjective findings of infectious or inflammatory disorders of their sinonasal complex. In creating a surgical corridor for access to the skull base, six distinct anatomical locations within the sinonasal cavity were sampled by both swab and biopsy (Fig. 6, Table 4, Additional file 2: Table S15).

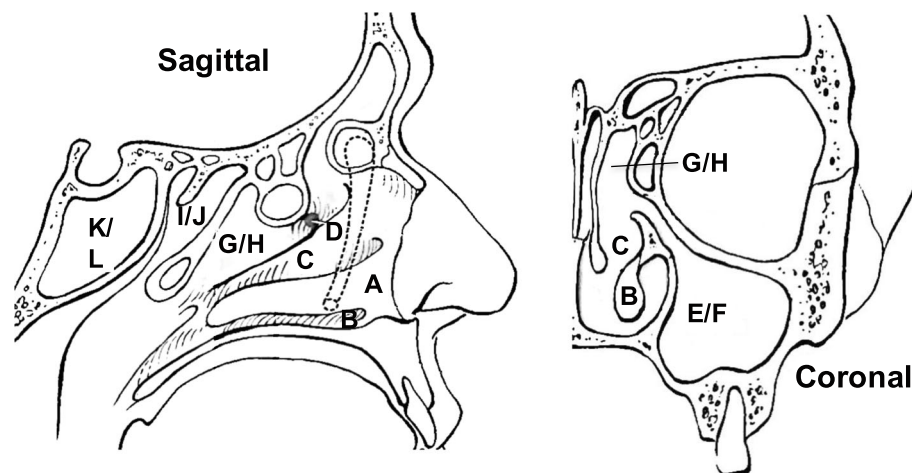


Fig. 6 Schematic diagrams in the sagittal and coronal planes of the human sinonasal cavity. Sites of sampling for microbiome analysis: deep nasal vestibule swab, deep to the vibrissae past the squamous mucosal epithelial junction (a), head of inferior turbinate swab (b), middle meatus swab (c), unciniate process biopsy (d), maxillary sinus swab (e) and biopsy (f), ethmoid sinus swab (g) and biopsy (h), superior meatus swab (i) and biopsy (j), and sphenoid sinus swab (k) and biopsy (l). Figure adapted from "Atlas of Endoscopic Sinus and Skull Base Surgery," ed. Palmer, J.N., Chiu, A.G., Adappa N.D. Elsevier, Philadelphia (2013)

To determine the bacterial constituents of the human sinonasal microbiome and the extent to which it varies among healthy individuals and among distinct sinonasal sites, we sequenced FL16S amplicons by PacBio from paired swabs and biopsies at the six anatomical sites from the 12 subjects (Additional file 2: Table S15, 122 specimens total across 12 individuals). No filtered CCS reads were generated from 23 samples, primarily from the maxillary sinus (both swabs and biopsies), suggesting little colonization of this site by bacteria (Fig. 6), and significantly fewer reads were collected from biopsy samples than from swab samples, potentially indicating lower overall bacterial load compared with the mucosal surface (Additional file 3: Figure S3). Filtering, clustering at an $EE \leq 1$, taxonomic assignment, and counts per OTU per sample were conducted as above (counting all filtered CCS reads against all non-chimera centroid OTU sequences via `usearch`). Complete information about counts per OTU per sample, as well as the taxonomic assignments of each centroid OTU, are in Additional file 2:

Table 4 Codes used for swab/biopsy sites

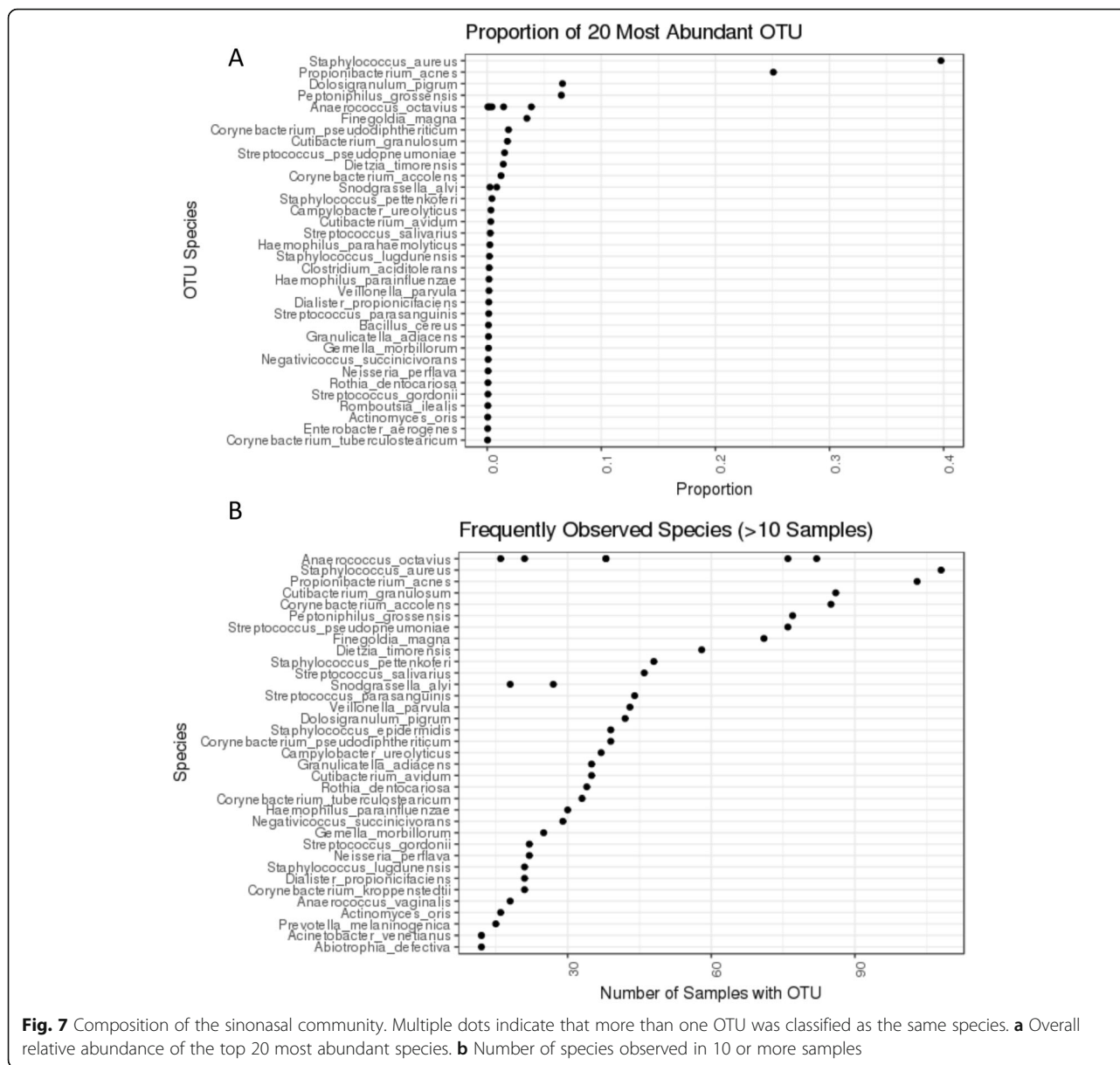
Swabs		Biopsies	
Code	Site	Code	Site
A	Nasal vestibule	B	Head of inferior turbinate tissue
C	Middle meatus	D	Uncinate process tissue
E	Maxillary sinus	F	Maxillary sinus tissue
G	Ethmoid culture (deep to ethmoid bulla)	H	Ethmoid tissue (deep to ethmoid bulla)
I	Superior meatus	J	Sphenoethmoidal recess tissue
K	Sphenoid	L	Sphenoid tissue

Tables S16–S18, and the results for all three communities have been incorporated into individual phyloseq objects in Additional file 5 (BEI), Additional file 6 (CAMI), and Additional file 7 (sino-nasal) [81].

The overall diversity of the sinonasal microbiomes collected here were relatively low. Across all specimens, clustering resulted in a total of 300 OTU (plus six centroids that were removed before classification by the CHIM2 filter), and the corresponding centroid OTU sequences were classified to 271 named species comprising 150 genera. Although 300 OTU were detected overall, the top 20 OTU comprised 96.7% of reads (Fig. 7), and only 61 OTU had >50 read counts summed across all >460K primer-match CCS reads. As previously seen, the dominant taxa in the sinonasal microbiome were *Staphylococcus* (OTU_2; see below) and *Propionibacterium acnes* (OTU_1), which together comprised 65.2% of all read counts [34]. Three of the top 20 OTU (and seven in total) were classified as *Anaerococcus octavius*, which suggests high variation among 16S rRNA genes within this species (Additional file 3: Figure S17).

We next investigated the relationship between species-level confidence values for each centroid and how many species are shared in the same dbOTU (Additional file 3: Figure S18). This analysis identified several OTU whose centroid assignment belonged to a dbOTU with only one or two species. These may represent other novel or poorly described species, or alternatively some may represent problems with the taxonomy assignments in the NCBI database.

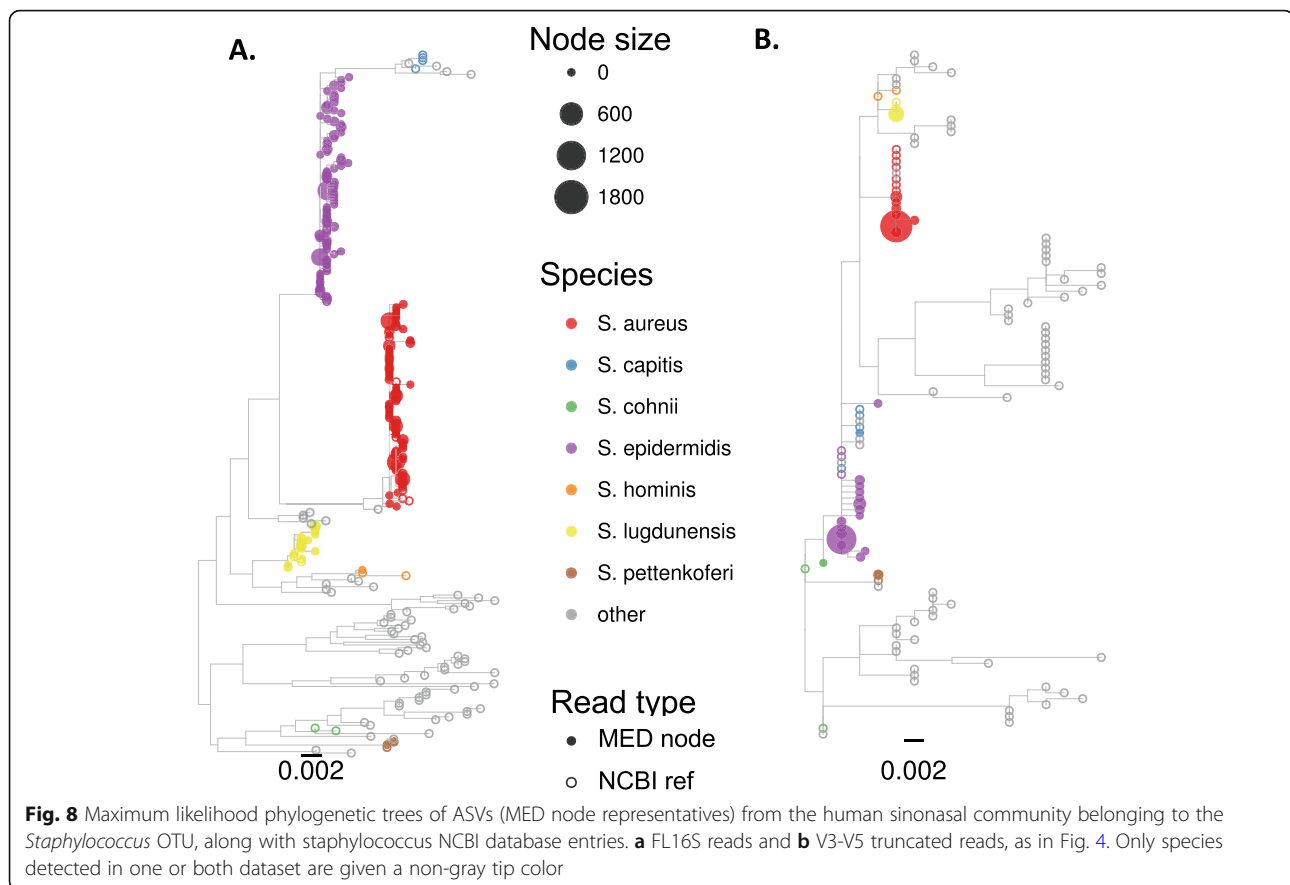
MED analysis of filtered CCS reads that had been assigned to the high abundance staphylococcal OTU



(whose centroid was assigned to *S. epidermidis*) further distinguished among distinct staphylococcal species within the human sinonasal samples, and this was improved when using FL16S compared to V3-V5 reads (Fig. 8). The presence within the sinonasal communities of additional close relatives to *S. aureus* and *S. epidermidis* clarified how V3-V5 truncated reads likely made some erroneous assignments, compared to FL16S. For example, examination of the trees in Fig. 8 suggests that the V3-V5 ASVs for *Staphylococcus capitus* and *Staphylococcus cohnii* are likely misclassified *S. epidermidis* sequences, and also that the *Staphylococcus hominis* ASV detected with FL16S reads was likely misclassified as *Staphylococcus lugdunensis* with the V3-V5 reads.

Variation in sinonasal microbial communities among subjects and anatomical sites

Considerable variation in microbial composition was seen among sinonasal specimens, ranging from three to 56 OTU per sample and from 50 to 108 per subject. The “core” sinonasal microbiome consisted of 11 OTU that were present across all 12 subjects, whereas most taxa were found in only a few individuals (Fig. 7b). For subsequent analyses of microbial diversity, absolute counts were normalized to relative abundances after first removing low yield samples and rare taxa, though results were qualitatively similar even with no filtering. We set a minimum sample size of 500 read counts, reducing the number of samples from 122 to 108 (the number of

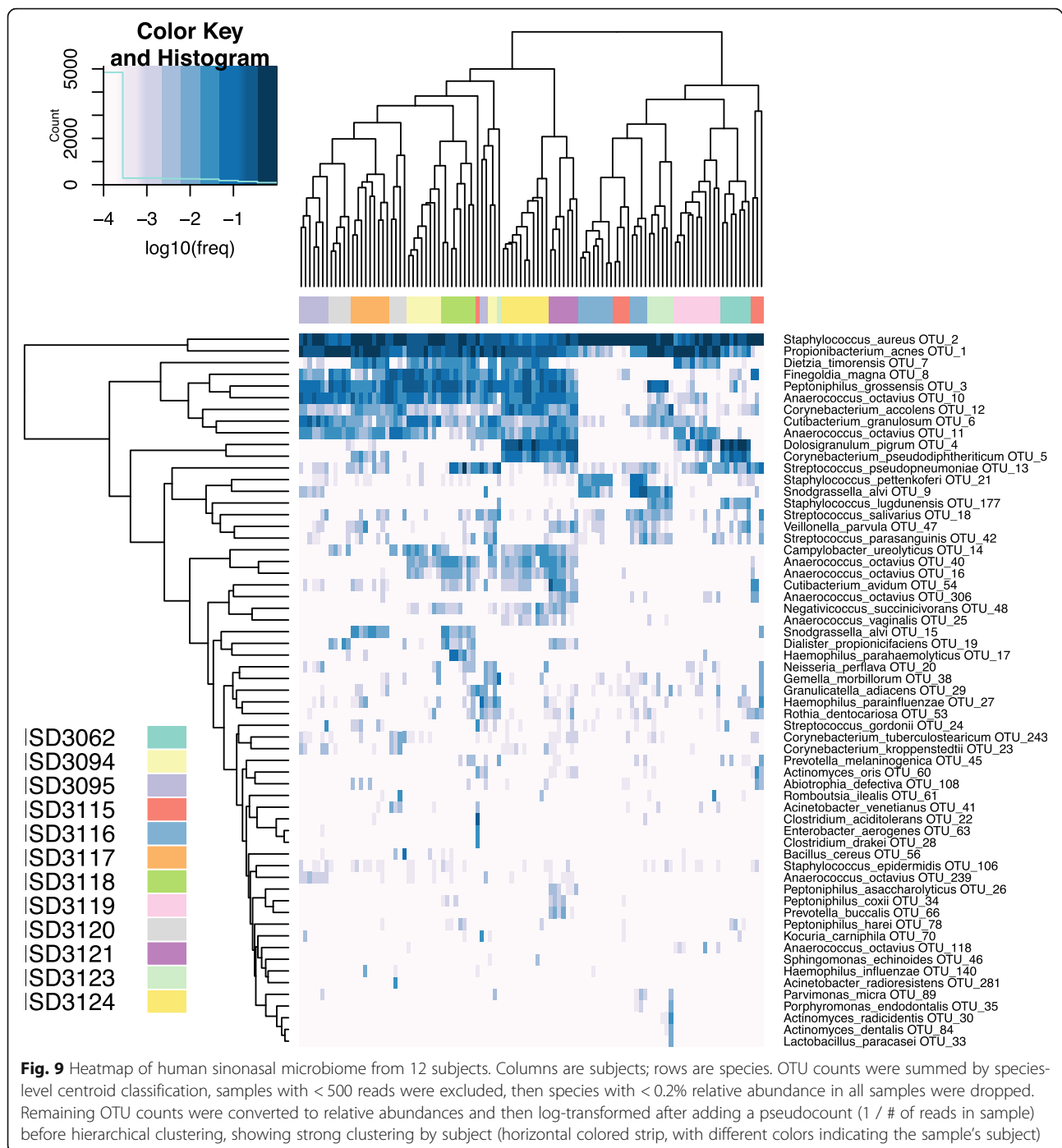


reads collected per sample was highly variable, ranging from 0 to 17,548, mean 3842 ± 3229). We also set a minimum OTU size of 50 read counts across the whole set of samples, reducing the number of taxa to only 59 OTU across the dataset. Though rare taxa may play important roles in the sinonasal microbiome, as has been shown in other environments [5], in the absence of dense longitudinal sampling, we could not tell whether these were resident to the sinonasal passages, transients, or contaminants.

Overall, the taxonomic profiles across samples were distinctly more similar within-subjects than within-site, as illustrated by hierarchical clustering and NMDS ordination of samples (Figs. 9 and 10a). This suggests that though the bacterial composition varies at distinct sub-anatomical sites, differences in microbial composition among individuals are much higher.

Because many OTU were only found in a subset of specimens, we next examined differences in the overall diversity of the samples with respect to subject, anatomical site, and whether obtained by swab or biopsy. Instead of using OTU richness (i.e., the total number of OTUs in each sample), we calculated Shannon's diversity index (which accounts for the relative abundance of distinct OTUs). Analysis-of-variance (ANOVA) of Shannon's

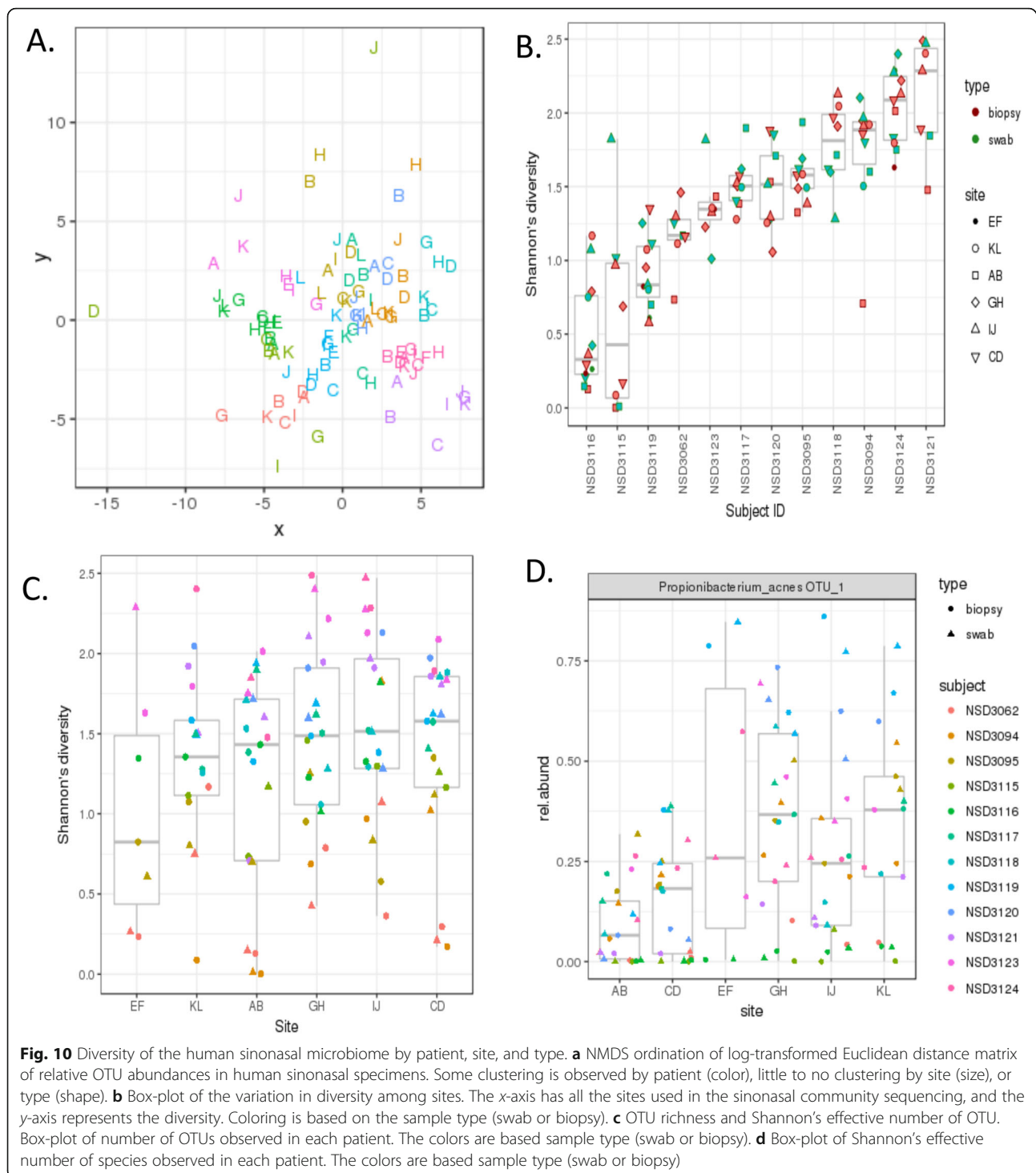
diversity found that, by far, the most important factor accounting for variation in Shannon's diversity was the subject the sample had come from (Fig. 10b, $p < 2e-16$). Sample type (swab versus biopsy) showed no significant effect ($p = 0.116$). Furthermore, although anatomical site was a significant contributor to the variance, no obvious trends were seen; variation among subjects was much higher (Fig. 10c, $p = 0.0082$). Swab/biopsy pairs from the same site and subject were extremely similar ($p = 0.96$), indicating no major shift in bacterial composition between the mucosal layer and the tissue immediately beneath, though the latter likely had fewer bacteria overall. These findings were robust to changing the filters used, to using genus- or species-level classifications, and to reformulating the ANOVA model with different factor orders and interaction terms. Furthermore, these results were not an artifact of undersampling in some samples, since there was no correlation between within-sample Shannon's diversity and sample read count (Additional file 3: Figure S19). A distinct test that accounts for under-sampled rare taxa may be more appropriate (via the breakaway package for R [82]), but due to the relatively low diversity of individual samples, we were unable to apply this test due to a requirement of seeing six consecutive frequency classes (not met



in any sample). Overall, these results suggest some underlying community structure in the sinonasal cavity, though much of this effect is hidden by the much larger differences in overall microbial composition among subjects.

Finally, to examine whether bacteria might partition differently within the sinonasal cavity, we performed ANOVA on transformed relative abundance measurements for each OTU. Only a single OTU showed a

significant effect by anatomical site (*Propionibacterium acnes*, $p = 0.009$ after Benjamini-Hochberg FDR correction), and none showed variation by swab versus biopsy. Interestingly, *P. acnes* was least abundant relative to other bacteria in the nasal vestibule (site AB) the largest and most aerated part of the sinonasal cavity, whereas its abundance often made up a major component of the bacterial signature at other less accessible sites (Fig. 10d).



Discussion

We report a novel high specificity pan-bacterial molecular diagnostic pipeline for profiling the bacterial composition of microbiome samples, applying amplification and sequencing of full-length 16S rRNA genes (FL16S) with the Pacific Biosciences (PacBio) platform. We exploit circular consensus sequencing (CCS), in which we obtain > 10

passes on average of each single molecule sequenced, resulting in CCS reads with exceptionally high quality. This single-molecule correction system is not possible on other modern DNA sequencers [18, 83]. Notably, our MCSMRT software is modular and can easily provide inputs to well-established and commonly used downstream microbiome analysis pipelines (namely QIIME and

Mothur) at several points before or after OTU clustering, ASV detection, taxonomic assignment, and abundance calculations.

Previous applications of PacBio to sequencing the 16S rRNA gene were initially hampered by higher error rates and insufficient polymerase processivity to leverage circular consensus sequencing [18, 19]. Subsequent improvements in PacBio sequencing chemistry have mostly overcome this [20, 21], and more recent efforts have shown the value of FL16S sequencing by PacBio for and identified the major considerations needed for handling PacBio instead of Illumina 16S reads [22–24]. This work extends and improves upon previous efforts in several ways:

- (1) We provide the flexible MCSMRT pipeline to handle processing, clustering, and taxonomic assignment of PacBio 16S reads after first identifying and implementing a series of stringent filters that eliminate many sources of sequencing artifacts. In particular we show that using only the highest fidelity consensus reads for OTU clustering (those with a cumulative expected error, EE, of ≤ 1) effectively eliminates over-calling the number of OTU, which has been a pervasive problem in methods using shorter partial 16S sequences [55], comparable to recent observations using PacBio sequencing of the same BEI mock community as we use here [23].
- (2) We generated new PacBio FL16S datasets for pipeline development and benchmarking, including monoculture controls from two species, low and high complexity mock communities, and 122 samples from the human sinonasal microbiome.
- (3) Because taxonomic assignments remain especially important in the study of human-associated bacteria, we developed a species-level taxonomic classifier for FL16S. To assign taxonomy and confidence values CCS reads, we created a custom-built database constructed from all available FL16S sequences at NCBI, since many commonly used 16S rRNA gene databases lack species-level classifications or lack FL16S genes for many taxa. This allowed us to use a uniform Linnaean hierarchy build a classifier that defined FL16S gene sequences to the species level along with associated confidence values. Our analysis showed improved accuracy and higher confidence when using FL16S sequences compared to partial sequences, which are typical when using short-read Illumina MiSeq 16S survey methods that normally capture only up to ~ 500 nt using paired-end sequencing (e.g., [62]).
- (4) We investigated the use of minimum entropy decomposition (MED) to detect amplicon sequence

variants (ASVs) as a way of distinguishing among closely related organisms [44, 84]. This found that decomposing OTUs into ASVs improved identification of closely related species, although the number of ASVs detected exceeded that expected within the BEI mock community. Some of this could be attributable to intragenomic variation among 16S rRNA gene copies, as seen with our *E. coli* monoculture positive controls, but we also suspect inflated ASV counts due to the particular error profile of CCS reads, which is still biased towards short indel variants, as well as the aggregated effect of true indel variation over the full length of the 16S rRNA gene when including diverse bacterial species. Future work to improve ASV detection from FL16S CCS reads via MED and/or DADA2 will likely require incorporating both pre-clustering, multiple alignment, and gap-handling during ASV detection.

We tested our experimental and bioinformatics pipeline on two distinct “mock communities,” showing that use of high-quality CCS reads from FL16S genes has exceptional precision and accuracy at identifying and quantifying the bacterial taxa in complex mixtures, which remains especially important in the study of clinically relevant bacteria. Most impressively, we correctly identified most species in the more complex CAMI community, making only three false identifications represented by only five CCS reads. Remarkably, we also accurately measured the relative abundance of most species in this complex CAMI community, indicating that our pipeline not only has high taxonomic specificity but also high accuracy for quantitative measures of species abundance in complex microbial communities. Underlining the accuracy and precision of our experimental and bioinformatics procedures, 99.81% of filtered reads from monoclonal positive control samples were correctly classified to the species level, and seven independent negative controls yielded almost no reads that passed our filters.

Following validation of our pipeline using complex mock communities, we applied MCSMRT to the human sinonasal microbiome, finding that the community has a relatively low complexity (with 61 OTU at a frequency of $> 0.1\%$ across samples); *Staphylococcus* species and *Propionibacterium acnes* dominated across subjects and anatomical sites [34]. Although microbial composition varied much more substantially among subjects than among anatomical sites in the same subject, we nevertheless observed trends in the overall diversity of different sites, with the easily accessible swabs just deep to the nasal vestibule overall reflecting the majority of the healthy sinonasal cavity with the least diverse and least

dominated by *P. acnes*. Importantly, we find that swab and biopsy sampling at the same site in the same subject have highly correlated microbial composition, indicating that invasive biopsy sampling is not needed. The large differences among the sinonasal microbiomes among healthy subjects will be of interest in future studies that examine links between sinonasal disease states (particularly chronic rhinosinusitis), bacterial composition, and the innate immune response [85–92]. Our results also show improved discrimination among closely related *Staphylococcus* species when using FL16S compared to the V3-V5 region alone.

Our results and others show that increasing the length of marker gene sequencing improves the taxonomic and phylogenetic resolution, and we expect that further improvements to sequence processing and analysis will greatly enhance studies that use ASV detection. Our use of primers targeting all nine variable regions was chosen to maximize specificity, but it may also have narrowed the overall breadth of bacterial diversity we could capture [93, 94], so future studies will investigate primer combinations that maximize both breadth and specificity. Sequencing on the PacBio platform remains more expensive than Illumina per read, though the price has dropped considerably since the introduction of the Sequel instrument and is expected to drop further when higher yield SMRTcells are released. Thus, we expect that the cost trade-off between higher specificity with PacBio (taxonomic and phylogenetic resolution) and higher sensitivity with Illumina (high yields) will rapidly diminish.

We further expect that CCS analysis can expand the scope of other marker-based taxonomic and phylogenetic studies, for example, through joint single-molecule sequencing of eukaryotic ITS and 18S rRNA genes from fungi to enrich and extend marker-based databases [95]. Beyond marker gene surveys, metagenomic shotgun sequencing efforts have shown the massive potential for simultaneous profiling of functional gene content and high resolution phylogenetic and taxonomic binning. Although these approaches often remain prohibitively expensive for profiling many host-associated microbiota and may be less amenable to use in clinical diagnostics, we note that metagenomic shotgun assembly and downstream analysis could potentially be greatly enhanced by use of high-quality CCS reads.

Methods

Ethics statement

Patients were recruited from the Division of Rhinology of the Department of Otorhinolaryngology—Head and Neck Surgery at the University of Pennsylvania with full approval of the Institutional Review Board (Protocol 800614). Informed consent was obtained during the pre-operative clinic visit or in the pre-operative waiting

room. Selection criteria for recruitment were patients undergoing sinonasal surgery for non-rhinologic disease entities, e.g., pituitary pathology or other cranial base pathologies.

Sinonasal sample collection

Sinonasal samples were obtained from patients undergoing sinonasal surgery for non-inflammatory and non-infectious indications (predominately pituitary tumors or other skull base neoplastic process) who had not received antibiotics in the preceding 8 weeks. The Institutional Review Board at The University of Pennsylvania School of Medicine provided full study approval, and informed consent was obtained pre-operatively from all patients. Sinonasal specimens were collected as both swabs (S) (BD ESwab collection and transport system) and tissue (T) (placed in MP lysing matrix tubes). Multiple locations in the sinonasal cavity were sampled including the nasal vestibule (S), inferior turbinate head (T), uncinate process (T), middle meatus (S), maxillary sinus (S)(T), ethmoid sinus (S)(T), superior meatus (S), superior turbinate (T), and sphenoid sinus (S)(T) for a maximum of 12 specimens per patient.

DNA extractions

Total DNA was isolated from all samples (swabs and biopsies) using DNeasy Blood & Tissue Kit (Qiagen) according to the manufacturers' recommendations with slight modifications. Biopsy material was incubated overnight at 56 °C with 570 µl ATL lysis buffer with 30 µl Proteinase K in a Lysing Matrix E tube (MP Biomedicals LLC), homogenized by SPEX 1600 MiniG (Fisher Sci.) for 10 min at 1500 Hz, and centrifuged 1 min × 13,000 rpm. Swab tubes were treated similarly but initially vortexed for 1 min and spun for 10 s and incubated for only 5 min at 56 °C prior to homogenization. DNA was eluted with 200 µl of the Elution Buffer. DNA quality and quantity were analyzed by agarose gel electrophoresis and Nanodrop 2000 spectrophotometry.

Control DNA samples

The BEI mock community was obtained through BEI Resources, NIAID, NIH as part of the Human Microbiome Project (<https://www.beiresources.org>): We used genomic DNA from Microbial Mock Community B (Even, Low Concentration), v5.1L, for 16S RNA Gene Sequencing, HM-782D. The complex CAMI mock community was obtained from the JGI, which had been constructed for the CAMI (Critical Assessment of Metagenomic Interpretation) Hosts Community Challenge for Assessing Metagenomes. Human DNA was isolated from the U937 lymphoblast lung cell line as an off-target control template.

FL16S rDNA PCR reactions

Amplifications were performed using 1 µl total DNA as template, universal 16S primers F27 and R1492 with four sets of asymmetric barcodes at 0.25 µM (Additional file 2: Table S3) [96, 97], and GoTaq Hot Start Master Mix (Promega) or AccuPrime Taq High Fidelity Polymerase with 1 µl of 10 mM dNTP Mix (Fisher Sci.) in 50 µl final volume. Cycling conditions were 94 °C, 3 min; then 22 or 35 cycles of 94 °C 30 s, 54 °C 30 s, 72 °C 2 min; following by 5 min final elongation at 72 °C. PCR products were cleaned with AxyPrep™ MagPCR (Corning Life Sciences) according to the manufacturer's protocol and eluted in 40 µl of water. Cleaned PCR products were quantified using both Quant-iT™ dsDNA Assay Kit, high sensitivity (Invitrogen) on BioTek™ FLx800™ Microplate Fluorescence Reader, and AccuClear Ultra High Sensitivity sDNA Quantitation Kit (Biotium). Based on the results, amplicons were normalized to the same concentration prior to pooling amplicons with distinct barcodes into multiplexed sets of two to four samples per pool.

Pacific Biosciences circular consensus sequencing

Library construction used Pacific Biosciences (PacBio) SMRTbell™ Template Prep Kit V1 on normalized pooled PCR products, and sequencing was performed using the PacBio RS II platform using protocol "Procedure & Checklist - 2 kb Template Preparation and Sequencing" (part number 001-143-835-06). DNA Polymerase Binding Kit P6 V2 was used for sequencing primer annealing and polymerase binding. SMRTbell libraries were loaded on SMRTcells V3 at final concentration 0.0125 nM using the MagBead kit. DNA Sequencing Reagent V4 was used for sequencing on the PacBio RS II instrument, which included MagBead loading and stage start. Movie times were 3 h for all SMRTcells. PacBio sequencing runs were set up using RS Remote PacBio software and monitored using RS Dashboard software. Sequencing performance and basic statistics were collected using SMRT® Analysis Server v2.3.0.

Pre-clustering pipeline

MCSMRT accepts CCS data from the PacBio RSII sequencer and is divided into pre-clustering and clustering steps (Fig. 1, Pre-clustering Pipeline, Additional file 1). Sequences were generated using the reads of insert (ROI) protocol within Pacific Biosciences SMRT® Analysis Server; reads which had four or fewer CCS passes were removed. To further filter low quality or off-target sequences, reads failing three filters were removed: (a) CCS reads outside the range of 500 to 2000 bp, (b) those that aligned to the hg19 human genome with bwa v0.7.10-r789, and (c) those that did not match both primer sequences with usearch v8.1.1861 [14, 52]. Primer sequences were then trimmed, and reads were oriented

5' to 3' with respect to 16S rRNA transcription. Venn diagrams defining read filtration subsets were created using Venny [98].

OTU clustering, taxonomic classification of centroids, and OTU abundances

OTUs were generated using the UPARSE [40] algorithm in the usearch software, using parameters tuned for full-length sequence. In short, reads were de-replicated (and the number, or size, of identical sequences tracked in the header), then sorted by abundance. OTUs were iteratively created at a threshold of 3% divergence from any other existing OTU centroid sequence (i.e., reads within 97% similarity to an existing OTU centroid became a member of an existing cluster; otherwise, a new OTU was formed with that sequence) (De novo OTU Clustering Pipeline, Additional file 1).

To obtain a database capable of providing a species-level classification of the full-length sequences, all sequences annotated as FL16S genes were downloaded from NCBI in October 2015, and taxonomies were inferred from each read's 16S gid identifier via the associated txid. This newly formatted database contained species-level taxonomic information for OTU classification (16S rRNA Microbial Database, Additional file 1). Representative OTU sequences were assigned a taxonomy using a utax classifier built from this database.

Chimeric sequences were removed during the clustering process based on previously seen OTU centroid sequences (CHIM1 filtering), followed by removal of chimeric centroid OTU using uchime to filter the final OTU sequences using the RDP "gold" sequences [53].

OTU abundance was determined using usearch for filtered reads prior to the expected error threshold, reported as CCS read counts assigned to each centroid OTU. 16S rRNA copy number for the BEI community was estimated from provided quality control data, and OTU abundance by 16S rRNA copy number was calculated in R (Additional file 2: Table S1, Fig. 8). Expected OTU abundances for the CAMI datasets used the rrndb database to obtain a predicted 16S rRNA gene copy number for each taxon, using the mean of values at the lowest matching taxonomic level.

Sub-OTU and phylogenetic methods

A 3% divergence cutoff for OTU clustering is commonly used in comparing various partial 16S fragments [5, 9, 99, 100]. To further examine how individual reads were related to one another, mafft [101] alignments of individual genus/species sequences (including sequences from both OTU and matching NCBI database) were created. Pairwise distance matrices from the alignments were created using ape v4.1 [102] and seqinr v 3.3-6 [103]. Neighbor-joining trees were created using vegan v2.4-0

and ggtree v1.8.1. Data were visualized using ggplot2 v2.2.1 [104–106]. Additionally, maximum likelihood trees (FastTree v2.1.8) were calculated and visualized with ggtree [107]. OTU alignments were further used for ASV analysis using MED v2.1 decomposition, as described (Identification of Amplicon Sequence Variants by Minimum Entropy Decomposition, Phylogenetic inference, Additional file 1).

Ecological analyses of the healthy sinonasal microbiome

Measures of ecological diversity (number of species observed, Shannon's diversity index) were calculated for each sample using vegan before and after filtering to eliminate samples with < 500 CCS reads and OTU with < 50 CCS reads across all samples. Count tables were transformed to relative abundances prior to calculating dissimilarity and distance matrices by either Bray-Curtis or Euclidean distance metrics. NMDS ordinations of samples were generated with vegan, and heatmaps created using the gplots package for R.

Data availability

MCSMRT and tutorial data is available at <https://github.com/jpearl01/mcsmrt>. The curated NCBI FL16S database is available at <https://drive.google.com/file/d/1UaWvDnVf-GOOtL3ld4B0t5v7H5igB0To/view>. All FL16S sequencing data was deposited at NCBI under BioProject PRJNA 369775; Biosample IDs are in Additional file 2: Tables S4 and S8). The finished assembly of our lab stock of *E. coli* K12 MG1655 was deposited under the same BioProject, accession ID CP032667. For the new datasets collected here, OTU tables, sample Info, and phylogenetic trees are found as phyloseq objects in Additional files 5, 6, and 7. Additional V3-V5 sequence reads from Illumina MiSeq for the same BEI mock community was acquired from [61], found under NCBI BioProject PRJNA242424. Overlapping paired-end reads were joined with COPE, changing default parameters to allow for longer overlap (up to 250 base pairs) [108]. Reads were then imported into MCSMRT and run through the default pipeline. Additional FL16S PacBio CCS reads were acquired from [24], available from the JGI Genome Portal under accession "Phylotag." Shotgun assemblies of CAMI community members were obtained from the Joint Genome Institute Genome Portal. QIIME and Mothur pipelines were performed on FL16S, V3-V5 regions datasets as described (Mothur OTU Clustering Pipeline, QIIME OTU Clustering Pipeline, Additional file 1). CAMI reference genomes were acquired from the CAMI website in July of 2016. FL16S gene prediction was accomplished using rnammer v1.2 on all assembled contigs [109].

Additional files

Additional file 1: Supplementary Text. (DOCX 56 kb)

Additional file 2: **Table S1.** Expected composition of "Even" DNA Mock Community from BEI. **Table S2.** Expected composition of DNA Mock Community from JGI-CAMI. **Table S3.** PCR Primers Sequence. **Table S4.** Filtering stats for mock communities. **Table S5.** FL16S Denovo OTU Counts from QIIME2 with Greengenes v13_8 Classification. **Table S6.** FL16S Closed OTU Counts from QIIME2 with Greengenes v13_8 Classification. **Table S7.** BEI v3-v5 16S Denovo OTU Counts from QIIME2 with Greengenes v13_8 Classification. **Table S8.** BEI v3v5 16S Closed OTU Counts from QIIME2 with Greengenes v13_8 Classification. **Table S9.** OTU Counts and Taxonomic Classification for BEI Mock Community Classified With Mothur v1.35. **Table S10.** OTU Count and Taxonomic Classification on Mock Community from Singer et al. **Table S11.** Error Analysis of Reads Mapping to Positive Control *E. coli* MG1655 Reference. **Table S12.** MCSMRT Classification of CAMI Mock Community. **Table S13.** Classification Accuracies of Cami Community for Full-Length and v3-v5 Truncated 16S. **Table S14.** Monophy results for multi-species genera in the CAMI community. **Table S15.** Sinonasal dataset with patient samples collected from sites as shown in Table 4. **Table S16.** Sinonasal Sample Statistics and NCBI BioSampleID. **Table S17.** Sinonasal OUT Taxonomic Classification and Confidence. **Table S18.** Read counts mapping to each OTU centroid for all sinonasal samples. (XLSX 443 kb)

Additional file 3: **Figure S1.** Effect of primary filters on the number of reads. **Figure S2.** Insert size distribution. **Figure S3.** Effects of Host DNA on bacterial 16S yield. **Figure S4.** Total CCS yield vs. DNA yield vs. PCR yields. **Figure S5.** CCS reads mapping to the human genome from the sinonasal samples. **Figure S6.** Primer matching filters. **Figure S7.** Primer matching truncation and nucleotide variability against positive control *E. coli* forward and reverse primer matches. **Figure S8.** Primer matching truncation and nucleotide variability against positive control *A. tumefaciens* forward and reverse primer matches. **Figure S9.** Histogram of the number of species per DB cluster in NCBI. **Figure S10.** The effect of PCR cycle number and polymerase choice on OTU abundances in the BEI mock community. **Figure S11.** The effect of PCR cycles and polymerase on abundances of chimeric molecules in the BEI mock community. **Figure S12.** Substitution errors in BEI mock community. **Figure S13.** *E. coli* MG1655 16S copies and MED analysis. **Figure S14.** Phylogenetic trees of well-resolved multi-species genera *Clostridium* and *Desulfovibrio*. **Figure S15.** Phylogenetic trees of poorly-resolved multi-species genera *Azotobacter* and *Nonlabens*. **Figure S16.** Phylogenetic trees of multi-species genera with improved species resolution using FL16S for *Algoriphagus* and *Salegentibacter*. **Figure S17.** Phylogeny of *Anaerococcus*MED nodes from the sinonasal communities plus NCBI database entries. **Figure S18.** Relationship between species-level confidence in centroid assignments and the number of species in the matching dbOTU. **Figure S19.** Effective Number of Species as a function of sample read depth. (PDF 3167 kb)

Additional file 4: Results of clustering the curated NCBI database with Single, Average, and Complete clustering methods. (XLSX 1953 kb)

Additional file 5: BEI phyloseq object. (RDS 2 kb)

Additional file 6: CAMI phyloseq object. (RDS 14 kb)

Additional file 7: Sinonasal phyloseq object. (RDS 32 kb)

Abbreviations

BEI: Biological and Emerging Infections Resources Program; CAMI: Critical Assessment of Metagenome Interpretation; CCS: Circular consensus sequence; EE: Expected error; JGI: Joint Genome Institute; MCSMRT: Microbiome Classifier using Single Molecule Real-time Sequencing; NCBI: National Center for Biotechnology Information; NIAID: National Institute of Allergy and Infectious Diseases; NMDS: Non-metric multidimensional scaling; nt: Nucleotide; OTU: Operational taxonomic unit; Pacbio: Pacific Biosciences; RDP: Ribosomal Database Project; ROI: Reads of insert; rrrDB: Ribosomal RNA Operon Copy Number Database

Acknowledgements

We thank Danielle Reed for useful motivating discussions. We also thank Jason Limbo for technical support, Katherine Kercher for purified genomic

DNA from U937 cell line lymphoblast lung from human cells, and Carol Hope for help with preparation of the manuscript. We also thank Robert Edgar for responsive help with the UPARSE pipeline. The BEI Bacterial Mock Community was obtained through BEI Resources, NIAID, NIH as part of the Human Microbiome Project. Very special thanks to Julio Corral Jr. and Axel Visel at JGI Production Facility for providing the CAMI mock community.

Funding

Funding for this work was provided by R01DC013588 to NAC, R01DC02148, U01DK082316, and NCI HHSN261201600383P to GDE.

Availability of data and materials

The dataset(s) supporting the conclusions of this article is(are) available in the NCBI SRA repository, [unique persistent identifier and hyperlink to dataset(s) in <http://> format].

Software

Project name: MCSMRT

Project home page: <https://github.com/jpearl01/mcsmrt>

Archived version: DOI or unique identifier of archived software or code in repository (e.g. enodo).

Operating system(s): Linux 64 bit

Programming language: ruby

Other requirements: usearch v.8.1.1861

License: GPL

Any restrictions to use by non-academics: none

For databases, this section should state the web/ftp address at which the database is available and any restrictions to its use by non-academics.

Authors' contributions

JPE, NAC, GDE, and JCM conceived and designed study. JPE, AB, and RLE coded the software. NDA, JNP, NAC, and GDE provided the materials and reagents. JPE, JCM, AB, GDE, NAC, JK, SB, RLE, and JH Wrote the manuscript. All authors carried out the experiments and analyses. All authors read and approved the final manuscript.

Ethics approval and consent to participate

The Institutional Review Board at The University of Pennsylvania School of Medicine provided full study approval, and informed consent was obtained pre-operatively from all patients.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Author details

¹Department of Microbiology & Immunology, Centers for Genomic Sciences and Advanced Microbial Processing, Drexel University College of Medicine, 245 N 15th Street, Philadelphia, PA 19102, USA. ²Veteran's Administration Medical Center, Department of Otorhinolaryngology: Head and Neck Surgery, University of Pennsylvania Perelman School of Medicine, 3400 Spruce Street, 5 Ravdin, Philadelphia, PA 19104-4283, USA.

Received: 13 February 2018 Accepted: 2 October 2018

Published online: 23 October 2018

References

- Woese CR. Bacterial evolution. *Microbiol Rev.* 1987;51(2):221–71.
- Woese CR, Fox GE. Phylogenetic structure of the prokaryotic domain: the primary kingdoms. *Proc Natl Acad Sci U S A.* 1977;74(11):5088–90.
- Lane DJ, Pace B, Olsen GJ, Stahl DA, Sogin ML, Pace NR. Rapid determination of 16S ribosomal RNA sequences for phylogenetic analyses. *Proc Natl Acad Sci U S A.* 1985;82(20):6955–9.
- Olsen GJ, Woese CR. Ribosomal RNA: a key to phylogeny. *FASEB J.* 1993;7(1):113–23.
- Welch DBM, Mark Welch DB, Huse SM. Microbial diversity in the deep sea and the underexplored “rare biosphere”. In: *Handbook of molecular microbial ecology II*; 2011. p. 243–52.
- Weisburg WG, Barns SM, Pelletier DA, Lane DJ. 16S ribosomal DNA amplification for phylogenetic study. *J Bacteriol.* 1991;173(2):697–703.
- Terlizzi A, Anderson MJ, Bevilacqua S, Fraschetti S, Włodarska-Kowalczuk M, Ellingsen KE. Beta diversity and taxonomic sufficiency: do higher-level taxa reflect heterogeneity in species composition? *Divers Distrib.* 2009;15(3):450–8.
- Anderson MJ, Ellingsen KE, McArdle BH. Multivariate dispersion as a measure of beta diversity. *Ecol Lett.* 2006;9(6):683–93.
- Schloss PD. The effects of alignment quality, distance calculation method, sequence filtering, and region on the analysis of 16S rRNA gene-based studies. *PLoS Comput Biol.* 2010;6(7):e1000844.
- Barriuso J, Valverde JR, Mellado RP. Estimation of bacterial diversity using next generation sequencing of 16S rDNA: a comparison of different workflows. *BMC Bioinformatics.* 2011;12:473.
- Kim M, Morrison M, Yu Z. Evaluation of different partial 16S rRNA gene sequence regions for phylogenetic analysis of microbiomes. *J Microbiol Methods.* 2011;84(1):81–7.
- Di Bella JM, Bao Y, Gloor GB, Burton JP, Reid G. High throughput sequencing methods and analysis for microbiome research. *J Microbiol Methods.* 2013;95(3):401–14.
- Liu Z, DeSantis TZ, Andersen GL, Knight R. Accurate taxonomy assignments from 16S rRNA sequences produced by highly parallel pyrosequencers. *Nucleic Acids Res.* 2008;36(18):e120.
- Youssef N, Sheik CS, Krumholz LR, Najjar FZ, Roe BA, Elshahed MS. Comparison of species richness estimates obtained using nearly complete fragments and simulated pyrosequencing-generated fragments in 16S rRNA gene-based environmental surveys. *Appl Environ Microbiol.* 2009;75(16):5227–36.
- Gilbert JA, Meyer F, Antonopoulos D, Balaji P, Brown CT, Brown CT, Desai N, Eisen JA, Evers D, Field D, et al. Meeting report: the terabase metagenomics workshop and the vision of an Earth microbiome project. *Stand Genomic Sci.* 2010;3(3):243–8.
- Teeling H, Glöckner FO. Current opportunities and challenges in microbial metagenome analysis—a bioinformatic perspective. *Brief Bioinform.* 2012;13(6):728–42.
- Wang J, Jia H. Metagenome-wide association studies: fine-mining the microbiome. *Nat Rev Microbiol.* 2016;14(8):508–22.
- Fichot EB, Norman RS. Microbial phylogenetic profiling with the Pacific Biosciences sequencing platform. *Microbiome.* 2013;1(1):10.
- Mosher JJ, Bernberg EL, Shevchenko O, Kan J, Kaplan LA. Efficacy of a 3rd generation high-throughput sequencing platform for analyses of 16S rRNA genes from environmental samples. *J Microbiol Methods.* 2013;95(2):175–81.
- Babauta JT, Atci E, Ha PT, Lindemann SR, Ewing T, Call DR, Fredrickson JK, Beyenal H. Localized electron transfer rates and microelectrode-based enrichment of microbial communities within a phototrophic microbial mat. *Front Microbiol.* 2014;5:11.
- Mosher JJ, Bowman B, Bernberg EL, Shevchenko O, Kan J, Korlach J, Kaplan LA. Improved performance of the PacBio SMRT technology for 16S rDNA sequencing. *J Microbiol Methods.* 2014;104:59–60.
- Wagner J, Coupland P, Browne HP, Lawley TD, Francis SC, Parkhill J. Evaluation of PacBio sequencing for full-length bacterial 16S rRNA gene classification. *BMC Microbiol.* 2016;16(1):274.
- Schloss PD, Jenior ML, Koumpouras CC, Westcott SL, Highlander SK. Sequencing 16S rRNA gene fragments using the PacBio SMRT DNA sequencing system. *PeerJ.* 2016;4:e1869.
- Singer E, Bushnell B, Coleman-Derr D, Bowman B, Bowers RM, Levy A, Gies EA, Cheng JF, Copeland A, Klenk HP, et al. High-resolution phylogenetic microbial community profiling. *ISME J.* 2016;10(8):2020–32.
- D'Amore R, Ijaz UZ, Schirmer M, Kenny JG, Gregory R, Darby AC, Shakya M, Podar M, Quince C, Hall N. A comprehensive benchmarking study of protocols and sequencing platforms for 16S rRNA community profiling. *BMC Genomics.* 2016;17:55.
- Abreu NA, Nagalingam NA, Song Y, Roediger FC, Pletcher SD, Goldberg AN, Lynch SV. Sinus microbiome diversity depletion and *Corynebacterium tuberculoostearicum* enrichment mediates rhinosinusitis. *Sci Transl Med.* 2012;4(151):151ra124.
- Aurora R, Chatterjee D, Hentzleman J, Prasad G, Sindwani R, Sanford T. Contrasting the microbiomes from healthy volunteers and patients with

- chronic rhinosinusitis. *JAMA Otolaryngol Head Neck Surg.* 2013;139(12):1328–38.
28. Stephenson MF, Mfuna L, Dowd SE, Wolcott RD, Barbeau J, Poisson M, James G, Desrosiers M. Molecular characterization of the polymicrobial flora in chronic rhinosinusitis. *J Otolaryngol Head Neck Surg.* 2010;39(2):182–7.
 29. Stressmann FA, Rogers GB, Chan SW, Howarth PH, Harries PG, Bruce KD, Salib RJ. Characterization of bacterial community diversity in chronic rhinosinusitis infections using novel culture-independent techniques. *Am J Rhinol Allergy.* 2011;25(4):e133–40.
 30. Boase S, Foreman A, Cleland E, Tan L, Melton-Kreft R, Pant H, Hu FZ, Ehrlich GD, Wormald PJ. The microbiome of chronic rhinosinusitis: culture, molecular diagnostics and biofilm detection. *BMC Infect Dis.* 2013;13:210.
 31. Lee JT, Frank DN, Ramakrishnan V. Microbiome of the paranasal sinuses: update and literature review. *Am J Rhinol Allergy.* 2016;30(1):3–16.
 32. Bezerra TF, Padua FG, Gebrim EM, Saldiva PH, Voegels RL. Biofilms in chronic rhinosinusitis with nasal polyps. *Otolaryngol Head Neck Surg.* 2011;144(4):612–6.
 33. Feazel LM, Robertson CE, Ramakrishnan VR, Frank DN. Microbiome complexity and *Staphylococcus aureus* in chronic rhinosinusitis. *Laryngoscope.* 2012;122(2):467–72.
 34. Ramakrishnan VR, Feazel LM, Gitomer SA, Ir D, Robertson CE, Frank DN. The microbiome of the middle meatus in healthy adults. *PLoS One.* 2013;8(12):e85507.
 35. Biswas K, Hoggard M, Jain R, Taylor MW, Douglas RG. The nasal microbiota in health and disease: variation within and between subjects. *Front Microbiol.* 2015;9:134.
 36. Bassis CM, Tang AL, Young VB, Pynnonen MA. The nasal cavity microbiota of healthy adults. *Microbiome.* 2014;2:27.
 37. Paju S, Bernstein JM, Haase EM, Scannapieco FA. Molecular analysis of bacterial flora associated with chronically inflamed maxillary sinuses. *J Med Microbiol.* 2003;52(Pt 7):591–7.
 38. Power DA, Burton JP, Chilcott CN, Tagg JR, Dawes PJ. Non-culture-based analysis of bacterial populations from patients with chronic rhinosinusitis. *J Clin Microbiol.* 2005;43(11):5822–4.
 39. Kaspar U, Kriegeskorte A, Schubert T, Peters G, Rudack C, Pieper DH, Wos-Oxley M, Becker K. The culturome of the human nose habitats reveals individual bacterial fingerprint patterns. *Environ Microbiol.* 2016;18(7):2130–42.
 40. Edgar RC. UPARSE: highly accurate OTU sequences from microbial amplicon reads. *Nat Methods.* 2013;10(10):996–8.
 41. Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski RA, Oakley BB, Parks DH, Robinson CJ, et al. Introducing mothur: open-source, platform-independent, community-supported software for describing and comparing microbial communities. *Appl Environ Microbiol.* 2009;75(23):7537–41.
 42. Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, Fierer N, Peña AG, Goodrich JK, Gordon JI, et al. QIIME allows analysis of high-throughput community sequencing data. *Nat Methods.* 2010;7(5):335–6.
 43. Eren AM, Morrison HG, Lescault PJ, Reveillaud J, Vineis JH, Sogin ML. Minimum entropy decomposition: unsupervised oligotyping for sensitive partitioning of high-throughput marker gene sequences. *ISME J.* 2015; 9(4):968–79.
 44. Eren AM, Maignien L, Sul WJ, Murphy LG, Grim SL, Morrison HG, Sogin ML. Oligotyping: differentiating between closely related microbial taxa using 16S rRNA gene data. *Methods Ecol Evol.* 2013;4(12):1111–19.
 45. Szczyrba A, Hofmann P, Belmann P, Koslicki D, Janssen S, Droge J, Gregor I, Majda S, Fiedler J, Dahms E, et al. Critical assessment of metagenome interpretation—a benchmark of metagenomics software. *Nat Methods.* 2017; 14(11):1063–71.
 46. Stackebrandt E, Goebel BM. Taxonomic note: a place for DNA-DNA reassociation and 16S rRNA sequence analysis in the present species definition in bacteriology. *Int J Syst Evol Microbiol.* 1994;44(4):846–9.
 47. Konstantinidis KT, Tiedje JM. Genomic insights that advance the species definition for prokaryotes. *Proc Natl Acad Sci U S A.* 2005;102(7):2567–72.
 48. Nguyen N-P, Warnow T, Pop M, White B. A perspective on 16S rRNA operational taxonomic unit clustering using sequence similarity. *NPJ Biofilms Microbiomes.* 2016;2:16004.
 49. Edgar RC. Updating the 97% identity threshold for 16S ribosomal RNA OTUs. *Bioinformatics.* 2018;34(14):2371–75.
 50. Kim M, Oh H-S, Park S-C, Chun J. Towards a taxonomic coherence between average nucleotide identity and 16S rRNA gene sequence similarity for species demarcation of prokaryotes. *Int J Syst Evol Microbiol.* 2014;64(Pt 2): 346–51.
 51. Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJ, Holmes SP. DADA2: high-resolution sample inference from Illumina amplicon data. *Nat Methods.* 2016;13(7):581–3.
 52. Edgar RC. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics.* 2010;26(19):2460–1.
 53. Edgar RC, Haas BJ, Clemente JC, Quince C, Knight R. UCHIME improves sensitivity and speed of chimera detection. *Bioinformatics.* 2011;27(16): 2194–200.
 54. Cole JR, Wang Q, Fish JA, Chai B, McGarrell DM, Sun Y, Brown CT, Porras-Alfaro A, Kuske CR, Tiedje JM. Ribosomal database project: data and tools for high throughput rRNA analysis. *Nucleic Acids Res.* 2014;42(Database issue):D633–42.
 55. Schloss PD, Westcott SL. Assessing and improving methods used in operational taxonomic unit-based approaches for 16S rRNA gene sequence analysis. *Appl Environ Microbiol.* 2011;77(10):3219–26.
 56. Mysara M, Vandamme P, Props R, Kerckhof FM, Leys N, Boon N, Raes J, Monsieurs P. Reconciliation between operational taxonomic units and species boundaries. *FEMS Microbiol Ecol.* 2017;93(4). <https://doi.org/10.1093/femsec/fix029>.
 57. Edgar RC, Flyvbjerg H. Error filtering, pair assembly and error correction for next-generation sequencing reads. *Bioinformatics.* 2015;31(21):3476–82.
 58. 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.* 2010;12(1):118–23.
 59. Gloor GB, Hummelen R, Macklaim JM, Dickson RJ, Fernandes AD, MacPhee R, Reid G. Microbiome profiling by Illumina sequencing of combinatorial sequence-tagged PCR products. *PLoS One.* 2010;5(10):e15406.
 60. Bokulich NA, Subramanian S, Faith JJ, Gevers D, Gordon JI, Knight R, Mills DA, Caporaso JG. Quality-filtering vastly improves diversity estimates from Illumina amplicon sequencing. *Nat Methods.* 2013;10(1):57–9.
 61. Salipante SJ, Kawashima T, Rosenthal C, Hoogestraat DR, Cummings LA, Sengupta DJ, Harkins TT, Cookson BT, Hoffman NG. Performance comparison of Illumina and ion torrent next-generation sequencing platforms for 16S rRNA-based bacterial community profiling. *Appl Environ Microbiol.* 2014;80(24):7583–91.
 62. Kozich JJ, Westcott SL, Baxter NT, Highlander SK, Schloss PD. Development of a dual-index sequencing strategy and curation pipeline for analyzing amplicon sequence data on the MiSeq Illumina sequencing platform. *Appl Environ Microbiol.* 2013;79(17):5112–20.
 63. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, Peplies J, Glöckner FO. The SILVA ribosomal RNA gene database project: improved data processing and web-based tools. *Nucleic Acids Res.* 2013;41(Database issue):D590–6.
 64. Yilmaz P, Parfrey LW, Yarza P, Gerken J, Pruesse E, Quast C, Schweer T, Peplies J, Ludwig W, Glöckner FO. The SILVA and “all-species living tree project (LTP)” taxonomic frameworks. *Nucleic Acids Res.* 2014;42(Database issue):D643–8.
 65. Perraudeau F, Dudoit S, Bullard JH. Accurate determination of bacterial abundances in human metagenomes using full-length 16S sequencing reads: biorxiv.org; 2017. <https://doi.org/10.1101/228619>.
 66. Balvocite M, Huson DH. SILVA, RDP, Greengenes, NCBI and OTT - how do these taxonomies compare? *BMC Genomics.* 2017;18(Suppl 2):114.
 67. McDonald D, Price MN, Goodrich J, Nawrocki EP, DeSantis TZ, Probst A, Andersen GL, Knight R, Hugenholtz P. An improved Greengenes taxonomy with explicit ranks for ecological and evolutionary analyses of bacteria and archaea. *ISME J.* 2012;6(3):610–8.
 68. Edgar RC. Reference databases. 2015. Retrieved from <https://drive5.com/syntax/>.
 69. Lal D, Verma M, Lal R. Exploring internal features of 16S rRNA gene for identification of clinically relevant species of the genus *Streptococcus*. *Ann Clin Microbiol Antimicrob.* 2011;10:28.
 70. Jervis-Bardy J, Leong LE, Marri S, Smith RJ, Choo JM, Smith-Vaughan HC, Nosworthy E, Morris PS, O’Leary S, Rogers GB, et al. Deriving accurate microbiota profiles from human samples with low bacterial content through post-sequencing processing of Illumina MiSeq data. *Microbiome.* 2015;3:19.
 71. Fukushima M, Kakinuma K, Kawaguchi R. Phylogenetic analysis of *Salmonella*, *Shigella*, and *Escherichia coli* strains on the basis of the *gyrB* gene sequence. *J Clin Microbiol.* 2002;40(8):2779–85.

72. Roggenkamp A. Phylogenetic analysis of enteric species of the family Enterobacteriaceae using the *oriC*-locus. *Syst Appl Microbiol*. 2007;30(3):180–8.
73. Lluch J, Servant F, Paissé S, Valle C, Valière S, Kuchly C, Vilchez G, Donnadieu C, Courtney M, Burcelin R, et al. The characterization of novel tissue microbiota using an optimized 16S metagenomic sequencing pipeline. *PLoS One*. 2015; 10(11):e0142334.
74. Hang J, Desai V, Zavaljevski N, Yang Y, Lin X, Satya RV, Martinez LJ, Blaylock JM, Jarman RG, Thomas SJ, et al. 16S rRNA gene pyrosequencing of reference and clinical samples and investigation of the temperature stability of microbiome profiles. *Microbiome*. 2014;2:31.
75. Coenye T, Vandamme P. Intragenomic heterogeneity between multiple 16S ribosomal RNA operons in sequenced bacterial genomes. *FEMS Microbiol Lett*. 2003;228(1):45–9.
76. Sczyrba A, Hofmann P, Belmann P, Koslicki D, Janssen S, Dröge J, Gregor I, Majda S, Fiedler J, Dahms E, et al. Critical assessment of metagenome interpretation—a benchmark of metagenomics software. *Nat Methods*. 2017; 14(11):1063–71.
77. Stoddard SF, Smith BJ, Hein R, Roller BRK, Schmidt TM. *rrnDB*: improved tools for interpreting rRNA gene abundance in bacteria and archaea and a new foundation for future development. *Nucleic Acids Res*. 2015; 43(Database issue):D593–8.
78. Schwery O, O'Meara BC. *Monophy*: a simple R package to find and visualize monophyly issues. *PeerJ Computer Science*. 2016;2:e56.
79. Orlandi RR, Kingdom TT, Hwang PH, Smith TL, Alt JA, Baroody FM, Batra PS, Bernal-Sprekelsen M, Bhattacharyya N, Chandra RK, et al. International consensus statement on allergy and rhinology: rhinosinusitis. *Int Forum Allergy Rhinol*. 2016;6(Suppl 1):S22–209.
80. Smith SS, Evans CT, Tan BK, Chandra RK, Smith SB, Kern RC. National burden of antibiotic use for adult rhinosinusitis. *J Allergy Clin Immunol*. 2013;132(5):1230–2.
81. McMurdie PJ, Holmes S. *phyloseq*: an R package for reproducible interactive analysis and graphics of microbiome census data. *PLoS One*. 2013;8(4):e61217.
82. Willis A, Bunge J, Whitman T. Improved detection of changes in species richness in high diversity microbial communities. *J R Stat Soc Ser C*. 2017; 66(5):963–77.
83. Travers KJ, Chin CS, Rank DR, Eid JS, Turner SW. A flexible and efficient template format for circular consensus sequencing and SNP detection. *Nucleic Acids Res*. 2010;38(15):e159.
84. Mackey KRM, Hunter-Cevera K, Britten GL, Murphy LG, Sogin ML, Huber JA. Seasonal succession and spatial patterns of *Synechococcus* microdiversity in a salt marsh estuary revealed through 16S rRNA gene oligotyping. *Front Microbiol*. 2017;8:1496.
85. Lee RJ, Xiong G, Kofonow JM, Chen B, Lysenko A, Jiang P, Abraham V, Doghramji L, Adappa ND, Palmer JN, et al. T2R38 taste receptor polymorphisms underlie susceptibility to upper respiratory infection. *J Clin Invest*. 2012;122(11):4145–59.
86. Lee RJ, Kofonow JM, Rosen PL, Siebert AP, Chen B, Doghramji L, Xiong G, Adappa ND, Palmer JN, Kennedy DW, et al. Bitter and sweet taste receptors regulate human upper respiratory innate immunity. *J Clin Invest*. 2014; 124(3):1393–405.
87. Adappa ND, Farquhar D, Palmer JN, Kennedy DW, Doghramji L, Morris SA, Owens D, Mansfield C, Lysenko A, Lee RJ, et al. TAS2R38 genotype predicts surgical outcome in nonpolypoid chronic rhinosinusitis. *Int Forum Allergy Rhinol*. 2016;6(1):25–33.
88. Adappa ND, Zhang Z, Palmer JN, Kennedy DW, Doghramji L, Lysenko A, Reed DR, Scott T, Zhao NW, Owens D, et al. The bitter taste receptor T2R38 is an independent risk factor for chronic rhinosinusitis requiring sinus surgery. *Int Forum Allergy Rhinol*. 2014;4(1):3–7.
89. Lee RJ, Hariji BM, McMahan DB, Chen B, Doghramji L, Adappa ND, Palmer JN, Kennedy DW, Jiang P, Margolske RF, et al. Bacterial d-amino acids suppress sinonasal innate immunity through sweet taste receptors in solitary chemosensory cells. *Sci Signal*. 2017;10(495):eaam7703.
90. Mfuna Endam L, Filali-Mouhim A, Boisvert P, Boulet LP, Bosse Y, Desrosiers M. Genetic variations in taste receptors are associated with chronic rhinosinusitis: a replication study. *Int Forum Allergy Rhinol*. 2014;4(3):200–6.
91. Carey RM, Workman AD, Hatten KM, Siebert AP, Brooks SG, Chen B, Adappa ND, Palmer JN, Kennedy DW, Lee RJ, et al. Denatonium-induced sinonasal bacterial killing may play a role in chronic rhinosinusitis outcomes. *Int Forum Allergy Rhinol*. 2017;7(7):699–704.
92. Adappa ND, Truesdale CM, Workman AD, Doghramji L, Mansfield C, Kennedy DW, Palmer JN, Cowart BJ, Cohen NA. Correlation of T2R38 taste phenotype and in vitro biofilm formation from nonpolypoid chronic rhinosinusitis patients. *Int Forum Allergy Rhinol*. 2016;6(8):783–91.
93. Ong SH, Kukkilaya VU, Wilm A, Lay C, Ho EX, Low L, Hibberd ML, Nagarajan N. Species identification and profiling of complex microbial communities using shotgun Illumina sequencing of 16S rRNA amplicon sequences. *PLoS One*. 2013;8(4):e60811.
94. Klindworth A, Pruesse E, Schweer T, Peplies J, Quast C, Horn M, Glockner FO. Evaluation of general 16S ribosomal RNA gene PCR primers for classical and next-generation sequencing-based diversity studies. *Nucleic Acids Res*. 2013;41(1):e1.
95. Tedersoo L, Tooming-Klunderud A, Anslan S. PacBio metabarcoding of Fungi and other eukaryotes: errors, biases and perspectives. *New Phytol*. 2018;217(3):1370–85.
96. Frank JA, Reich CI, Sharma S, Weisbaum JS, Wilson BA, Olsen GJ. Critical evaluation of two primers commonly used for amplification of bacterial 16S rRNA genes. *Appl Environ Microbiol*. 2008;74(8):2461–70.
97. Stackebrandt E, Goodfellow M. *Nucleic acid techniques in bacterial systematics*. Hoboken: Wiley; 1991.
98. Jol SJ. Make a Venn Diagram. 2015. Retrieved from <https://www.stefanjl.nl/venny>.
99. Sun Y, Cai Y, Liu L, Yu F, Farrell ML, McKendree W, Farmerie W. ESPRIT: estimating species richness using large collections of 16S rRNA pyrosequences. *Nucleic Acids Res*. 2009;37(10):e76.
100. Huse SM, Mark Welch DB. Accuracy and quality of massively parallel DNA pyrosequencing. In: *Handbook of molecular microbial ecology I*; 2011. p. 149–55.
101. Katoh K, Misawa K, Kuma K-I, Miyata T. MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Res*. 2002;30(14):3059–66.
102. Paradis E, Claude J, Strimmer K. *APE: analyses of phylogenetics and evolution in R language*. *Bioinformatics*. 2004;20(2):289–90.
103. Charif D, Lobry JR. *SeqinR 1.0-2: a contributed package to the R project for statistical computing devoted to biological sequences retrieval and analysis*. In: *Biological and medical physics, biomedical engineering*; 2007. p. 207–32.
104. Yu G, Smith DK, Zhu H, Guan Y, Lam TT-Y. *ggtree: an R package for visualization and annotation of phylogenetic trees with their covariates and other associated data*. *Methods Ecol Evol*. 2016;8(1):28–36.
105. Whickam H, Sievert C. *ggplot2: Elegant Graphics for Data Analysis (Use R!) 2nd Edition*. New York: Singer; 2016.
106. Dixon P. *VEGAN, a package of R functions for community ecology*. *J Veg Sci*. 2003;14(6):927–30.
107. Price MN, Dehal PS, Arkin AP. *FastTree 2 – approximately maximum-likelihood trees for large alignments*. *PLoS One*. 2010;5(3):e9490.
108. Liu B, Yuan J, Yiu SM, Li Z, Xie Y, Chen Y, Shi Y, Zhang H, Li Y, Lam TW, et al. COPE: an accurate k-mer-based pair-end reads connection tool to facilitate genome assembly. *Bioinformatics*. 2012;28(22):2870–4.
109. Lagesen K, Hallin P, Rodland EA, Staerfeldt HH, Rognes T, Ussery DW. *RNAmmr: consistent and rapid annotation of ribosomal RNA genes*. *Nucleic Acids Res*. 2007;35(9):3100–8.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

