

RESEARCH

A framework for assessing 16S rRNA marker-gene survey data analysis methods using mixtures.

Nathan D. Olson^{1,2,3*}, M. Senthil Kumar^{2,3}, Shan Li⁴, Domenick J. Braccia^{2,3}, Stephanie Hao⁵, Winston Timp⁵, Marc L. Salit⁶, O. Colin Stine⁴ and Hector Corrada Bravo^{2,3,7}

*Correspondence: nolson@nist.gov

¹Biosystems and Biomaterials Division, National Institute of Standards and Technology, 100 Bureau Dr., Gaithersburg, Maryland, 20899 USA
Full list of author information is available at the end of the article

Abstract

Background: Analysis of 16S rRNA marker-gene surveys may be performed by a variety of bioinformatic pipelines and downstream analysis methods. However, appropriate assessment datasets and statistics are needed as there is limited guidance to decide between available analysis methods. Mixtures of environmental samples are useful for assessment as they provide values calculated from measurements of the unmixed samples and the mixture design that can be compared to values recovered by each bioinformatic method. While experiments mixing complex samples have been used to assess other sequencing methods such as RNAseq, they have yet to be used to assess 16S rRNA sequencing.

Results: We developed an assessment framework for 16S rRNA sequencing analysis methods based on a two-sample titration mixture dataset and metrics to evaluate OTU count table characteristics. Our qualitative assessment evaluates feature presence/absence exploiting features only present in unmixed samples or titrations by testing if random sampling can explain their observed relative abundance. Our quantitative assessment evaluates how well relative and differential abundance values agree with values expected from the mixture design. We evaluated count tables generated by three commonly used bioinformatic pipelines as demonstration: i) DADA2 a sequence inference method, ii) Mothur a *de novo* clustering method, and iii) QIIME which uses open-reference clustering. Qualitative assessment indicated that the majority of Mothur and QIIME features specific to unmixed samples or titrations were explained by random sampling alone but not DADA2 features. When combined with assessments of count table sparsity, these results indicate that DADA2 has a higher false negative rate whereas Mothur and QIIME have higher false positive rates. Quantitative assessment indicated that, overall, observed relative abundance and differential abundance values were consistent with expected values for all three pipelines. We also identified subsets of features measured with high error by all pipelines evaluated. We could not identify the source of bias in these poor performing features based on previously studied sources of bias, indicating that further analysis of potentially unknown and unaccounted for biases is warranted.

Conclusions: We developed a novel framework for assessing 16S rRNA marker-gene survey analysis methods based on mixture experiments. To demonstrate the assessment framework we evaluated count tables generated using three bioinformatic pipelines. The assessment framework developed for this study will serve as a valuable community resource for assessing 16S rRNA marker-gene survey bioinformatic methods.

Keywords: 16S rRNA gene; assessment; bioinformatic pipeline; normalization; differential abundance

Background

Targeted sequencing of the 16S rRNA gene is commonly used to characterize microbial communities. The 16S rRNA marker-gene-survey measurement process includes molecular steps to selectively target and sequence the 16S rRNA gene from prokaryotic organisms within a sample and computational steps [1] computational steps convert the raw sequence data into a count table of feature relative abundance values [1]. Both molecular and computational measurement processes contribute to the overall measurement bias and dispersion [2, 1, 3]. The need for datasets characterizing complex microbial communities with some degree “ground truth” has emerged in order to properly characterize the accuracy of the 16S rRNA marker-gene-survey measurement process.

Diverse bioinformatic pipelines used to generate count tables produce data with diverse characteristics. For example the commonly used QIIME, Mothur, and DADA2 pipelines produce feature sets and count tables with different characteristics. Mothur uses *de novo* clustering for feature inference [4, 5]. Pairwise distances used in clustering are calculated from a multiple sequence alignment. Quality filtered paired-end reads are merged into contigs, then aligned to a reference multiple sequence alignment, followed by the removal of uninformative positions. As a result the feature set representative sequences are shorter than the input amplicons. For the QIIME open-reference clustering pipeline merged paired-end reads are first assigned to reference cluster centers [6, 7]. Next, unassigned reads are clustered *de novo*. Unlike Mothur, the QIIME pipeline clustering method uses pairwise sequence distances calculated from pairwise sequence alignments. As a result, the QIIME pairwise distances are calculated using the full amplicon sequence, whereas Mothur pairwise distances are calculated using multiple sequence alignment with only informative positions. The DADA2 pipeline uses a probability model and maximization expectation algorithm for feature inference [8]. Unlike distance-based clustering methods employed by the Mothur and QIIME pipelines, DADA2 parameters determine if low abundance sequences are grouped with a higher abundance sequence.

Numerous studies have evaluated quantitative and qualitative characteristics of the 16S rRNA measurement process using mock communities, simulated data, and environmental samples. Mock communities are commonly used to assess the qualitative characteristics of the 16S rRNA sequencing measurement process [9]. The use of mock communities in this fashion shows that surveys often result in number of features that are significantly higher than the underlying features in the mock community [10]. The higher than expected number of features is often attributed to sequencing and PCR artifacts as well as reagent contaminants [3, 11]. A notable exception is count tables generated using feature inference methods, such as DADA2 [8]. Sequence inference methods which aim to reduce the number of features from sequence artifacts by using statistical models to group sequences by both similarity and abundance. Nonetheless, while mock communities are useful in this type of assessment, they lack the diversity and dynamic range of feature present in real samples [9].

Quantitative assessment of 16S rRNA sequence data using mock communities and simulated data is informative but provides an incomplete characterization of the measurement process. Results from relative abundance estimates using mock communities generated from mixtures of single organism’s DNA have shown taxonomic

specific effects where individual taxa are under or over represented in a sample.¹
 For example, Gram-negative bacteria have higher extraction efficiency compared to²
 Gram-positive bacteria, and are thus likely over represented in count tables[12, 13].³
 Mismatches in the primer binding sites are also responsible for taxonomic specific⁴
 biases [3, 14, 15]. Additionally, taxon specific biases due to sequence template prop-⁵
 erties such as GC content, secondary structure, and gene flanking regions have been⁶
 observed [16, 17, 15]. However, due to limited community complexity the applicabil-⁷
 ity of mock community assessment results to more complex environmental samples⁸
 is unknown. Environmental sample complexity can be modeled using simulated and⁹
 have been used to assess differential abundance methods, where specific taxa are¹⁰
 artificially over represented in one set of samples compared to another [18]. How-¹¹
 ever, using simulated data to assess log fold-change estimates only evaluates the¹²
 computational steps of the measurement process.¹³

Quantitative and qualitative assessment can also be performed using sequence¹⁴
 data generated from mixtures of environmental samples. While simulated data and¹⁵
 mock communities are useful in evaluating and benchmarking new methods, one¹⁶
 needs to consider that methods optimized for mock communities and simulated¹⁷
 data are not necessarily optimized for the sequencing error profile and feature di-¹⁸
 versity of real samples. Data from real environmental samples are often used to¹⁹
 benchmark new molecular laboratory and computational methods. However, with-²⁰
 out expected values for use in assessment, only measurement precision or agreement²¹
 with other methods can be evaluated. By mixing environmental samples, expected²²
 values are calculated using information from the unmixed samples and mixture²³
 design. Mixtures of environmental samples were previously used to evaluate gene²⁴
 expression measurements [19, 20, 21].²⁵

Here we present a framework for assessing computational methods used to analyze²⁶
 16S rRNA marker-gene-survey data. The framework is comprised of a 16S rRNA²⁷
 two-sample titration dataset, generated using mixtures of human stool sample DNA²⁸
 extracts, along with metrics to assess the quantitative and qualitative characteristics²⁹
 of count tables generated using marker-gene-survey computational methods. To³⁰
 demonstrate usage of this assessment framework, we evaluated three bioinformatic³¹
 pipelines. Both the dataset and metrics developed in this study are publicly available³²
 and can be used to evaluate and optimize new and existing bioinformatic pipelines.³³

Results³⁵

Assessment Framework³⁶

Our framework assesses the qualitative and quantitative characteristics of the 16S³⁷
 rRNA measurement process (Fig. 1). The framework evaluates count tables gener-³⁸
 ated by bioinformatic pipelines from a dataset developed specifically for use in this³⁹
 framework. The qualitative assessment provides insight into how much confidence⁴⁰
 a user can have in feature presence/absence. The quantitative assessment evaluates⁴¹
 the bias and variance of relative and differential abundance estimates.⁴²

Assessment Dataset - Mixture Design⁴⁴

Using mixtures of environmental samples we generated a dataset with expected⁴⁵
 values for use in our assessment framework. For mixture datasets, expected values⁴⁶

AssessmentFramework.pdf

Figure 1 Assessment Framework. A) Count tables evaluated by the assessment framework are generated from the assessment dataset using marker-gene survey bioinformatic pipelines. Count table rows are features identified by the bioinformatic pipeline and column are samples, four PCR replicates (labeled A-D) were sampled for PRE and POST and titrations, to simplify the diagram only three titrations are shown. B) Pictorial depiction of abundance values of the seven feature types observed and used in the assessment framework. C) Qualitative and quantitative assessment metrics used in the assessment framework. The artifactual feature proportion metric (AFP) is a qualitative assessment of feature presence/absence based on unmixed-specific or titration-specific artifactual features. Sparsity (SPAR) is a qualitative assessment of the proportion of observed features in each sample relative to the total observed features. Relative abundance metric (Rel) plot is a quantitative assessment of the relationship between the observed and expected relative abundance values. The difference is used to calculate the error rate $(|Obs - Exp|/Exp)$ from which the bias metric ($median(error)$) and variance metric ($RCOV$) are calculated. The differential abundance (Diff) metric assesses the relationship between the expected log fold-change and estimated log fold-change is shown. Points represent the log fold-change between two titrations, point text indicates the titrations compared. A linear model is fit to the data. The model fit information is used for the differential abundance bias ($1 - slope$) and variance metrics (R^2). Each feature type in (B) is labeled with the assessments shown in (C) in which they are employed.

experimentalDesign.pdf

Figure 2 Sample selection and experimental design for the two-sample titration 16S rRNA marker-gene-survey assessment dataset. A) Pre- and post-exposure (PRE and POST) samples from five vaccine trial participants were selected based on *Escherichia coli* abundance measured using qPCR and 454 16S rRNA sequencing (454-NGS), data from Pop et al. [22]. Counts represent normalized relative abundance values for 454-NGS and copies of the heat-labile toxin gene per μL , a marker gene for ETEC, for qPCR. PRE and POST samples are indicated with orange and green data points, respectively. Grey points are other samples from the vaccine trial time series. B) Proportion of DNA from PRE and POST samples in titration series samples. PRE samples were titrated into POST samples following a \log_2 dilution series. The NA titration factor represents the unmixed PRE sample. C) PRE and POST samples from the five vaccine trial participants, subjects, were used to generate independent two-sample titration series. The result was a total of 45 samples, 7 titrations + 2 unmixed samples times 5 subjects. Four replicate PCRs were performed for each of the 45 samples resulting in 190 PCRs.

can be obtained using information from unmixed samples and the mixture design. Our mixture dataset uses a two-sample titration mixture design, where DNA collected from five vaccine trial participants before and after exposure to pathogenic *Escherichia coli* was mixed following a \log_2 dilution series (Fig. 2). Each sample was sequenced in quadruplicate. For our two-sample titration mixture design, expected feature relative abundance is calculated using equation (1), where θ_i is the proportion of POST DNA in titration i , q_{ij} is the relative abundance of feature j in titration i , and the relative abundance of feature j in the unmixed PRE and POST samples is $q_{pre,j}$ and $q_{post,j}$. Throughout the rest of the manuscript, samples collected prior to and after *E. coli* exposure are referred to as PRE and POST respectively.

$$q_{ij} = \theta_i q_{post,j} + (1 - \theta_i) q_{pre,j} \quad (1)$$

Qualitative Assessment

The qualitative assessment shows how well pipelines differentiate true biological sequences from measurement process artifacts. Inadequate processing of artifacts results in false positive and false negative features where false positives are features

¹in a count table that are not present in the sequenced sample and false negative¹
²features are biological sequences in a sample not represented in the count table.²
³Our qualitative assessment methods characterize the artifactual feature proportion³
⁴(the frequency of artifactual features in a count table) by estimating the proportion⁴
⁵of *titration*- and *unmixed-specific* features (Fig. 1B) that cannot be explained by⁵
⁶sampling alone. We combine the artifactual feature proportion assessment results⁶
⁷with sparsity estimates to hypothesize whether the artifactual features are primarily⁷
⁸false positives or negatives. Sparsity is defined as the fraction of 0 valued cells in⁸
⁹the count table (Fig. 1C).⁹

¹¹*Quantitative Assessment*¹¹

¹²To evaluate count table abundance values, our quantitative assessment uses error,¹²
¹³bias, and variance metrics (Fig. 1C). Error metrics measure agreement between ob-¹³
¹⁴served and expected abundance values. The bias and variance metrics summarise¹⁴
¹⁵feature-level performance. Bias metrics summarise the overall agreement with ex-¹⁵
¹⁶pected values and the variance metric characterizes the distribution of the agree-¹⁶
¹⁷ment. Overall, pipeline performance is evaluated by comparing count table metric¹⁷
¹⁸distributions. Additionally, feature-level metrics are indicators of feature-specific¹⁸
¹⁹biases.¹⁹

²¹Assessment Dataset Characterization and Validation²¹

²²To assure the mixture dataset is suitable for use in our assessment framework, we²²
²³first validated the titration series and raw sequence data. The mixture dataset had²³
²⁴sufficient sample coverage, reads per sample, and read quality for use in our assess-²⁴
²⁵ment framework. The number of reads per sample and distribution of base quality²⁵
²⁶scores by position was consistent across subjects (Fig. S5). There were 8.9548×10^{12} ²⁶
²⁷(152,267 - 3,195) sequences per sample, median and range. Average base quality²⁷
²⁸score was greater than 30 over the length of the amplicon when considering both²⁸
²⁹forward and reverse reads (Fig. S5B).²⁹

³⁰ Additionally, we characterized subject specific differences to inform the interpre-³⁰
³¹tation of our assessment results. No subject specific differences in base quality score³¹
³²were observed (Fig. S5). However, average read depth was greater for E01JH004³²
³³compared to the other individuals (Fig. S5). Community composition differences³³
³⁴between PRE and POST samples and individuals was characterized using alpha³⁴
³⁵and beta diversity (Fig. S6). Overall alpha diversity was higher for POST except³⁵
³⁶for E01JH0011, though differences in diversity between PRE and POST varied by³⁶
³⁷individual. Based on the beta diversity the community composition within individ-³⁷
³⁸uals differed between the PRE and POST samples. Note that assessment metrics³⁸
³⁹defined above and results reported below are based on within subject comparisons.³⁹

⁴⁰ To validate the two-sample titration assessment dataset, we evaluated two as-⁴⁰
⁴¹sumptions about the titrations: (1) The samples were mixed volumetrically in a \log_2 ⁴¹
⁴²dilution series according to the mixture design. (2) The unmixed PRE and POST⁴²
⁴³samples have the same proportion of prokaryotic DNA. To validate the sample vol-⁴³
⁴⁴umetric mixing exogenous DNA (ERCC plasmids) were spiked into the unmixed⁴⁴
⁴⁵samples before mixing and quantified using qPCR (Fig. S1B). The stool samples⁴⁵
⁴⁶used to generate the mixtures have both eukaryotic (primarily human) DNA and⁴⁶

Table 1 Summary statistics for the different bioinformatic pipelines. DADA2 is a denoising sequence inference pipeline, QIIME is an open-reference clustering pipeline, and Mothur is a de-novo clustering pipeline. No template controls were excluded from summary statistics. Sparsity is the proportion of 0's in the count table. Features is the total number of OTUs (QIIME and Mothur) or SVs (DADA2) in the count. Sample coverage is the median and range (minimum-maximum) per sample total abundance. Drop-out rate is the proportion of reads removed while processing the sequencing data for each bioinformatic pipeline.

Pipelines	Features	Sparsity	Total Abundance	Drop-out Rate
DADA2	3144	0.93	68649 (1661-112058)	0.24 (0.18-0.59)
Mothur	38358	0.98	53775 (1265-87806)	0.4 (0.35-0.62)
QIIME	11385	0.94	25254 (517-46897)	0.7 (0.62-0.97)

prokaryotic DNA. If the proportion of prokaryotic DNA differs between the unmixed samples, then the amount of DNA from the unmixed samples in a titration targeted by 16S rRNA gene sequencing is not consistent with the mixture design. We quantified the proportion of prokaryotic DNA in the unmixed samples using a qPCR assay targeting the 16S rRNA gene (Fig. S1C).

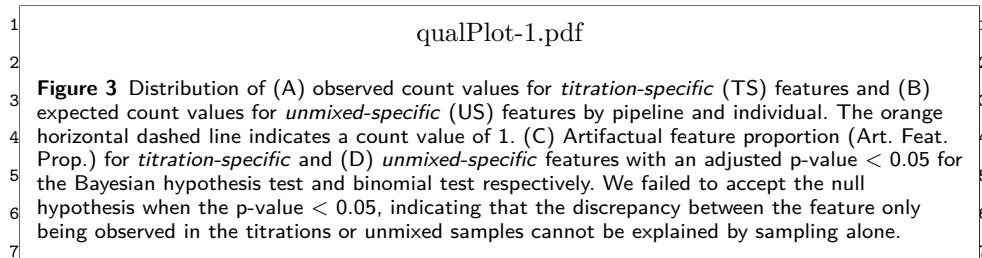
Our assessment dataset validation results indicated that the samples were volumetrically mixed according to the mixture design (Table S1) but prokaryotic DNA proportion varied across the titration series (Fig. S2). To account for deviations from the mixture design due to differences in the proportion of prokaryotic DNA in the unmixed samples, we estimated the proportion of POST in each titration using the 16S rRNA sequencing data (Fig. S3) and the estimated POST proportions were used in our assessment metric calculations. See Supplemental Material for the assessment dataset validation methods and results.

Count Table Assessment Demonstration

Next, we demonstrate the utility of our assessment framework on count tables generated using three different bioinformatic pipelines; DADA2, Mothur and QIIME. First, we provide high level summary statistics for initial insight into how the count tables differ. Next, we compare the assessment framework results for the three count tables.

Count Table Characteristics The count tables generated using the three bioinformatic pipelines vary in pre-processing and feature inference methods. These differences are reflected in the count table number of features, total abundance, and drop-out rate (Table 1, Fig. S7B). The pipelines evaluated employ different approaches for handling low quality reads resulting in large differences in the drop-out rate, fraction of raw sequences not included in the count table (Table 1). QIIME pipeline has the highest drop-out rate and number of features per sample but fewer total features than Mothur. The targeted amplicon region has a relatively small overlap region, 136 bp for 300 bp paired-end reads, compared to other commonly used amplicons [23, 24]. The high drop-out rate is due to low basecall accuracy at the ends of the reads especially the reverse reads resulting in a high proportion of unsuccessfully merged reads pairs (Fig. S5B). Further increasing the filter rate, QIIME excludes singletons (features only observed once in the dataset).

Feature taxonomic composition also varied by pipeline (Fig. S8). The three pipelines generated unique feature sets in terms of sequence length and amplicon position (see pipeline description). Therefore, we used feature taxonomic assignments



for cross-pipeline community composition comparison. Phylum and order relative abundance is similar across pipelines (Fig. S8A & B). The observed differences are attributed to different taxonomic classification methods and databases used by the pipelines. Regardless of the relative abundance threshold, most genera were unique to individual pipelines (Fig. S8C & D). Sets (shared taxa between pipelines) with QIIME had the fewest genera, excluding the DADA2-QIIME set. QIIME was the only pipeline to use open-reference clustering and the Greengenes database. Mothur and DADA2 both used the SILVA dataset. The Mothur and DADA2 pipeline use different implementations of the RDP naïve Bayesian classifier, which may be partially responsible for the Mothur, unclustered, and DADA2 differences.

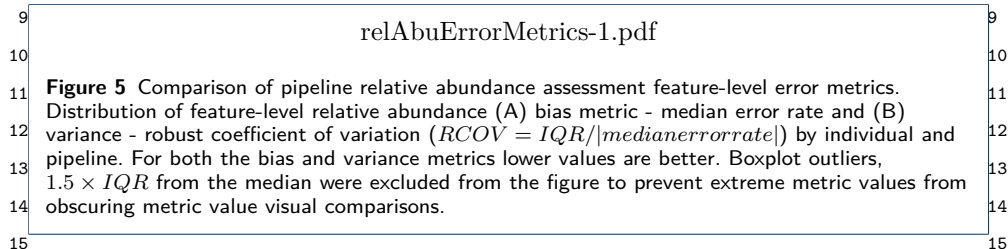
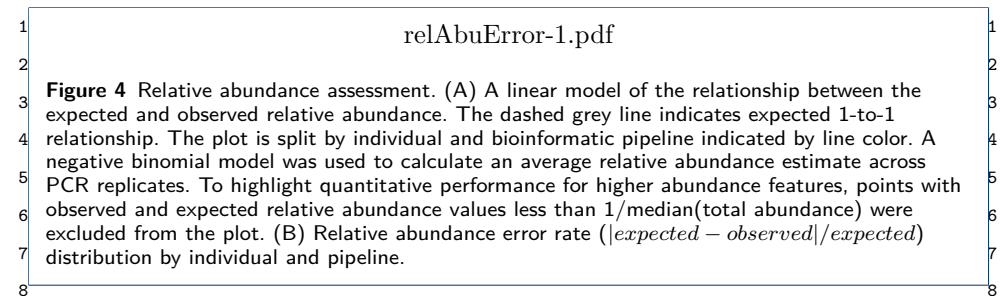
Qualitative Assessment

To evaluate feature presence-absence, the framework's qualitative assessment measures artifactual feature proportion and count table sparsity. Low abundance features present only in unmixed samples or titration samples are expected due to random sampling. *Unmixed-* and *titration-specific* features were observed for all pipelines (*titration-specific*: Fig. 3A, *unmixed-specific*: Fig. 3B). Overall, the DADA2 count table had the largest number of artifactual features (Table S3). A summary of the *titration-specific* artifactual features is provided in the supplementary material.

We next assessed the proportion of these artifactual features that could be explained by sampling effects alone. For our two-sample titration dataset, there were *unmixed-specific* features with expected counts not which could not be explained by sampling alone for all individuals and bioinformatic pipelines (Fig. 3C). However, the proportion of *unmixed-specific* features that could not be explained by sampling alone varied by bioinformatic pipeline. DADA2 had the highest proportion of *unmixed-specific* artifactual features whereas QIIME had the lowest proportion which is consistent with the distribution of *titration-specific* feature observed counts (Fig. 3D).

We expected this mixture dataset to be less sparse relative to other datasets due to the redundant nature of the samples where the 35 titration samples are derived directly from the 10 unmixed samples, along with four PCR replicates for each sample. We observed overall sparsity of 0.93 and 0.94 for DADA2 and QIIME respectively, and a higher value of 0.98 for Mothur 1).

To account for differences in microbial community composition across the five individuals we also measured sparsity at the individual level (Table S2). Sparsity at the individual-level is lower than overall sparsity for all three pipelines. In this case, average sparsity across individuals for 0.70 and 0.76 for DADA2 and Mothur, while QIIME had a lower average sparsity across individuals of 0.56. Differences in



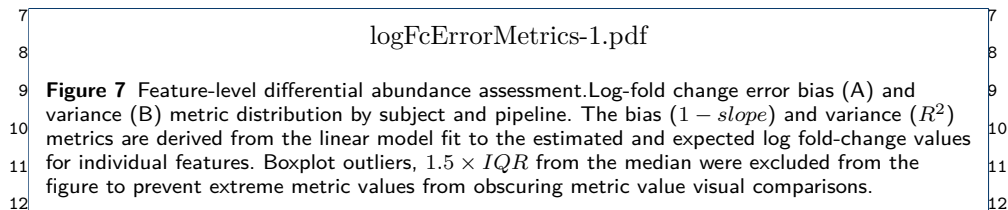
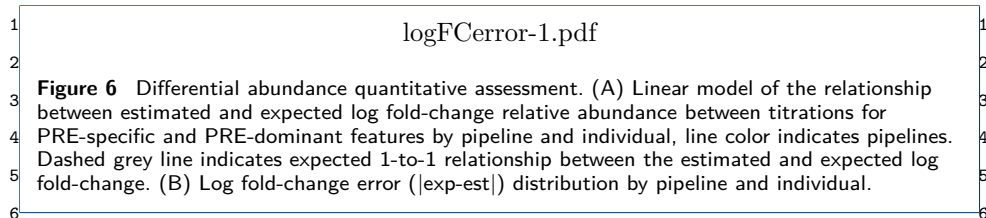
alpha and beta diversity for the five individual unmixed samples are consistent with individual level sparsity and therefore reflects differences in individual microbial community composition.

Based on the artifactual feature proportions and count table sparsity, DADA2 artifactual features are likely due to false negative features, whereas the Mothur and QIIME high sparsity values were attributed to false positive features. Based on the observed sparsity levels it is unlikely that any of the pipelines successfully filtered out a majority of the sequencing artifacts. Both unmixed- and titration-specific features that can and cannot be explained by sampling alone contribute to sparsity and the differences in the artifactual feature proportion and sparsity provide insight into how the pipelines treat sequencing artifacts.

Quantitative Assessment

Relative Abundance Assessment To assess count table feature relative abundance values, we evaluated the consistency of the observed and expected relative abundance estimates for a feature and titration as well as feature-level bias and variance. Only features observed in all PRE and POST PCR replicates and PRE and POST specific features were included in the analysis (Table S3). Overall, agreement between inferred and observed relative abundance was high for all individuals and bioinformatic pipelines (Fig. 4A). The error rate distribution was similarly consistent across pipelines, including long tails (Fig. 4B).

To assess quantitative accuracy across pipelines, we compared the feature-level relative abundance error rate bias and variance using mixed effects models. To control for subject specific differences, subject was included in the model as a random effect. Large bias and variance metric values were observed for all pipelines (Table S3). Feature-level relative abundance error rate bias (median error rate, Fig. 5A) was significantly different between pipeline, but no statistically significant differences were observed for the variance metric, ($RCOV = (IQR)/|\text{median}|$, Fig. 5B) across pipeline. The Mothur, DADA2, and QIIME feature-level biases were all significantly different from each other ($p < 1 \times 10^{-8}$). DADA2 had the lowest mean



feature-level bias (0.2), followed by Mothur (0.28), with QIIME having the highest bias (0.33) (5B). Large variance metric values were observed for all individuals and pipelines (Table S3). The feature-level variance was not significantly different between pipelines: Mothur = 0.83, QIIME = 0.71 and DADA2 = 1 (Fig. 5B).

Differential Abundance Assessment The agreement between log-fold change estimates and expected values were individual specific and consistent across pipelines (Fig. 6A). The individual specific effect was attributed to the fact that unlike relative abundance assessment, the inferred θ values were not used to calculate expected values. Inferred θ values were not used to calculate the expected values because all of the titrations and the θ estimates for the higher titrations were not monotonically decreasing. Using the inferred θ resulted in unrealistic expected log fold-change values, e.g., negative log-fold changes for PRE specific features. The log-fold change estimates and expected values were consistent across pipelines with one notable exception: for subject E01JH0011, the Mothur log fold-change estimates were more consistent with expected values than the other pipelines. However, as θ was not corrected for differences in the proportion of prokaryotic DNA between the un-mixed PRE and POST samples, it cannot be said whether Mothur's performance was better than the other pipelines.

The log fold-change error distribution was consistent across pipelines (Fig. 6B). There was a long tail of high error features in the error distribution for all pipelines and individuals. The log fold-change estimates responsible for the long tail could not be attributed to specific titration comparisons. Additionally, we compared error distributions for log-fold change estimates using different normalization methods. Error rate distributions, including the long tails, were consistent across normalization methods. Seeing as the long tail was observed for the unclustered data as well, the log-fold change estimates contributing to the long tail are likely due to a bias associated with the molecular portion of the measurement process and not the computational portion. Exploratory analysis of the relationship between the log fold-change estimates and expected values for individual features indicated that the long tails were attributed to feature specific performance.

Feature-level log fold-change bias and variance metrics were used to compare pipeline performance (Fig. 6). Similar to relative abundance, feature-level bias and

¹variance metrics are defined as the $1 - slope$ and R^2 for linear models of the estimated¹
²and expected log fold-change for individual features and all titration comparisons.²
³For the bias metric, $1 - slope$, the desired value is 0 (i.e., log fold-change estimate =³
⁴log fold-change expected), with negative values indicating the log-fold change was⁴
⁵consistently underestimated and positive values consistently overestimated. The⁵
⁶linear model R^2 value was used to characterize the feature-level log fold-change⁶
⁷variance as it indicates consistency between log fold-change estimates and expected⁷
⁸values across titration comparisons. To compare bias and variance metrics across⁸
⁹pipelines, mixed-effects models were used. The log fold-change bias and variance⁹
¹⁰metrics were not significantly different between pipelines (Bias: $F = 0, 2.51, p = 10$
¹¹ $0.99, 0.08, 6B, \text{ Variance: } F = 47.39, 0.23, p = 0, 0.8, \text{ Fig. } 6C$).¹¹

¹³Discussion¹³

¹⁴Mixtures of environmental samples have been used to assess RNAseq and microarray¹⁴
¹⁵gene expression measurements [19, 20, 21]. However, this is the first time mixtures¹⁵
¹⁶have been used to assess microbiome measurement methods. We developed a novel¹⁶
¹⁷assessment framework utilizing a mixture dataset for evaluating marker-gene-survey¹⁷
¹⁸computational methods (Fig. 1).¹⁸

¹⁹Using mixtures of environmental samples, expected values for use in assessment¹⁹
²⁰can be obtained using information from unmixed samples and how the samples²⁰
²¹were mixed. Our assessment dataset follows a two-sample titration mixture design,²¹
²²where DNA collected from five vaccine trial participants before and after exposure²²
²³to pathogenic *Escherichia coli* was mixed following a \log_2 dilution series (Fig. 2).²³
²⁴Count table qualitative characteristics were assessed using relative abundance in-²⁴
²⁵formation for features observed only in titrations (titration-specific) and unmixed²⁵
²⁶samples (unmixed-specific) (Fig. 1B). Statistical tests were used to determine if the²⁶
²⁷absence of unmixed-specific features from titrations or absence of titration-specific²⁷
²⁸features from unmixed samples could be explained by random sampling. Count ta-²⁸
²⁹bles were quantitatively assessed by comparing observed feature relative abundance²⁹
³⁰and feature differential abundance estimates to expected values. Quantitative per-³⁰
³¹formance was characterized using error rate, along with feature-level bias variance³¹
³²metrics we developed (Fig. 1C).³²

³⁴Count Table Assessment Demonstration³⁴

³⁵We demonstrated our assessment framework on count tables generated by three³⁵
³⁶commonly used bioinformatic pipelines, QIIME, Mothur, and DADA2. The objec-³⁶
³⁷tive of any pipeline is to differentiate true biological sequences from measurement³⁷
³⁸process artifacts along with accurate abundance estimates. Our qualitative assess-³⁸
³⁹ment results, when combined with sparsity information provides a new method for³⁹
⁴⁰evaluating how well bioinformatic pipelines account for sequencing artifacts without⁴⁰
⁴¹loss of true biological sequences. Additionally, our quantitative assessment results⁴¹
⁴²identified previously unknown feature specific biases in abundance estimates.⁴²

⁴³The qualitative assessment evaluates if titration- and unmixed-specific features⁴³
⁴⁴can be explained by random sampling alone (Fig. 1B). Titration- and unmixed-⁴⁴
⁴⁵specific features not explained by sampling are artifacts of the measurement pro-⁴⁵
⁴⁶cess. These artifacts can be viewed as false-positives, not representative of actual⁴⁶

sequences in a sample, or false-negatives, actual sequences in a sample not represented by count table features. Artifacts can be PCR errors such as chimeras, reads with high sequencing error rates, or cross sample contamination [25, 26, 27]. Count table sparsity information (the proportion of zero-valued cells) provides additional insight into the qualitative assessment results.

A high false negative rate provides an explanation for DADA2's high proportion of artifact titration- and unmixed-specific features and count table having comparable sparsity to the other pipelines despite having significantly fewer features (Fig. S5 and Table 1). The DADA2 feature inference algorithm may be aggressively grouping lower abundance true sequences with higher abundance sequences. As a result, the low abundance sequences are not present in samples leading to increased sparsity and high abundance unmixed- and titration-specific features. This aggressive grouping of sequences is a design choice made by the algorithm developers. The DADA2 documentation states that the default setting for `OMEGA_A` is conservative to prevent false positives at the cost of increasing false negatives [8]. Using the qualitative assessment methods described here, a user can adjust the `OMEGA_A` parameter to obtain a false-negative rate appropriate for their study.

While the relative abundance bias metric was significantly different between pipelines, overall, pipeline choice had minimal impact on the quantitative assessment results when accounting for subject-specific deviations in the proportion of prokaryotic DNA from PRE and POST samples in a titration from the mixture design. Outlier features (those with extreme bias and variance metrics) were observed for all pipelines and both abundance assessments.

Outlier features could not be attributed to bioinformatic pipelines and are likely due to biases in the molecular biology part of the measurement process. Outlier features are unlikely pipeline artifacts as they were observed in count tables generated using the unclustered pipeline as well as standard bioinformatic pipelines. Additionally, we were unable to attribute outlier features to relative abundance values, log fold-change between unmixed samples, and sequence GC content. Furthermore, features with extreme metric values were not limited to any specific taxonomic group or phylogenetic clade. PCR amplification bias (a well-known source of bias in the molecular biology part of the measurement process) is one possible explanation for the outlier features [28]. Mismatches in the primer binding regions impact PCR efficiency and are a potential cause for poor feature-specific performance [29]. Additional research is needed before outlier features can be attributed to mismatches in the primer binding regions.

Based on our assessment results, we suggest using DADA2 for feature-level abundance analysis, e.g. differential abundance testing. While DADA2 performed poorly in our qualitative assessment, the pipeline performed better in the quantitative assessment compared to the other pipelines. Additionally, the DADA2 poor qualitative assessment results due to false-negative features are unlikely to negatively impact feature-level abundance analysis. When determining which pipeline to use for a study, users should consider whether minimizing false positives (DADA2) or false negatives (Mothur) is more appropriate for their study objectives. Based on our findings we find that users of DADA2 can be more confident that an observed feature represents a member of the microbial community and not a measurement

artifact, but careful examination of sequences assigned to features of interest should still be performed.

Using Mixtures to Assess 16S rRNA Sequencing - Lessons Learned

There are limitations using our assessment dataset, these include: (1) Lack of agreement between the proportion of prokaryotic DNA from the unmixed samples in the titrations and the mixture design. (2) The mixture design resulted in a limited number of features and range of expected log-fold changes. These limitations are described below along with recommendations for addressing them in future studies.

Differences in the proportion of prokaryotic DNA in the samples used to generate the two-sample titrations series resulted in differences between the true mixture proportions and mixture design. We attempted to account for differences in mixture proportion from mixture design by using sequence data to estimate mixture proportions similar to how mRNA proportions in RNA samples were used in a previous mixture study [19]. We used an assay targeting the 16S rRNA gene to detect changes in the concentration of prokaryotic DNA across titrations, but were unable to quantify the proportion of prokaryotic DNA in the unmixed samples using qPCR data. Using the 16S rRNA sequencing data, we inferred the proportion of prokaryotic DNA from the POST sample in each titration. However, the uncertainty and accuracy of the inference method are not known, resulting in an unaccounted for source of error.

A better method for quantifying sample prokaryotic DNA proportion or using samples with consistent proportions would increase confidence in the expected value and, in-turn, error metric accuracy. Limitations in the prokaryotic DNA qPCR assay's concentration precision limits the assay's suitability for use in mixture studies. Digital PCR provides a more precise alternative to qPCR and is, therefore, a more appropriate method. Alternatively using samples where the majority of the DNA is prokaryotic would minimize this issue. Mixtures of environmental samples can also be used to assess shotgun metagenomic methods as well. As shotgun metagenomics is not a targeted approach, differences in the proportion of prokaryotic DNA in a sample would not impact the assessment results in the same way as 16S rRNA marker-gene-surveys.

Using samples from a vaccine trial allowed for the use of a specific marker with an expected response, *E. coli*, during methods development. However, the high level of similarity between the PRE and POST unmixed samples resulted in a limited number of features that could be used in the quantitative assessment results. Using more diverse samples to generate mixtures would address this issue. Alternatively, instead of mixing PRE and POST samples from the same individual, mixing PRE and POST samples from different individuals would have resulted in additional features for use in our quantitative assessment. While unmixed sample similarity impacts the number of features that can be used in the quantitative assessment, the qualitative assessment is not impacted by unmixed sample similarity. Finally, a symmetric mixture design, for example one with unmixed PRE and POST ratios of 1:4, 1:2, 1:1, 2:1, and 4:1, would provide a larger dynamic range of abundance values for assessing both PRE and POST specific features.

1Conclusions

2Our assessment framework can be used to evaluate and characterize 16S rRNA
3marker-gene survey analysis methods, in particular count tables produced by any
416S rRNA bioinformatic pipeline. We demonstrated our assessment framework with
5three commonly used bioinformatic pipelines. Our qualitative assessment results in-
6dicated that the QIIME and Mothur pipelines produced count table with more false-
7positive features whereas the DADA2 count table had more false-negative features.
8Overall the three pipelines performed well in our quantitative assessment. How-
9ever, feature-level analysis identified poorly performing features and the sources of
10bias responsible for this poor feature-level quantitative performance are unknown.
11Therefore, feature-level results for any 16S rRNA marker-gene survey should be
12interpreted with care. Addressing both of these issues requires advances in both the
13molecular biology and computational components of the measurement process.

15Methods

16Assessment Framework

17To assess the qualitative and quantitative performance of marker-gene survey
18analysis methods we developed a framework utilizing our two-sample titration
19dataset (Fig. 1). Qualitative assessment evaluates feature presence-absence. The
20quantitative assessment evaluates the relative and differential abundance estimates.

22Assessment Dataset - Mixture Design

23To provide a dataset with real-world complexity and expected values for qualitative
24and quantitative assessment we used mixtures of environmental samples. Samples
25collected at multiple timepoints during a Enterotoxigenic *E. coli* (ETEC) vaccine
26trial [30] were used to generate a two-sample titration dataset (Fig. 2). Samples
27from five trial participants were selected for our two-sample titration dataset. Trial
28participants (subjects) and sampling timepoints were selected based on *E. coli* abun-
29dance data collected using qPCR and 16S rRNA sequencing from Pop et al. [22].
30Only individuals with no *E. coli* detected in samples collected from trial partici-
31pants prior to ETEC exposure (PRE) were used for our two-samples titrations. Post
32ETEC exposure (POST) samples were identified as the timepoint after exposure
33to ETEC with the highest *E. coli* concentration for each subject (Fig. 2A). Due to
34limited sample availability, for E01JH0016 the timepoint with the second highest
35*E. coli* concentration was used as the POST sample. Independent titration series
36were generated for each subject. POST samples were titrated into PRE samples
37with POST proportions of 1/2, 1/4, 1/8, 1/16, 1/32, 1/1,024, and 1/32,768 (Fig.
382B). Unmixed (PRE and POST) sample DNA concentration was measured using
39NanoDrop ND-1000 (Thermo Fisher Scientific Inc. Waltham, MA USA). Unmixed
40samples were diluted to 12.5 ng/ μ L in tris-EDTA buffer before mixing. The result-
41ing titration series was composed of 45 samples, seven titrations and two unmixed
42samples for each of the five subjects.

43The 45 samples were processed using the Illumina 16S library protocol (16S
44Metagenomic Sequencing Library Preparation, posted date 11/27/2013, down-
45loaded from <https://support.illumina.com>). This protocol specifies an initial
46PCR of the 16S rRNA gene, followed by a sample indexing PCR, sample concen-
47tration normalization, and sequencing.

¹ A total of 192 16S rRNA PCR assays were sequenced across two 96-well plates¹
²including four PCR replicates per sample and 12 no-template controls. The ini-²
³tial PCR assay targeted the V3-V5 region of the 16S rRNA gene, Bakt_341F and³
⁴Bakt_806R [14]. The V3-V5 region is 464 base pairs (bp) long, with forward and⁴
⁵reverse reads overlapping by 136 bp, using 2 X 300 bp paired-end sequencing⁵
⁶[31] (<http://probase.csb.univie.ac.at>). Primer sequences include overhang⁶
⁷adapter sequences for library preparation (forward primer 5'- TCG TCG GCA⁷
⁸GCG TCA GAT GTG TAT AAG AGA CAG CCT ACG GGN GGC WGC AG -⁸
⁹3' and reverse primer 5'- GTC TCG TGG GCT CGG AGA TGT GTA TAA GAG⁹
¹⁰ACA GGA CTA CHV GGG TAT CTA ATC C - 3'). Kapa HiFi HotStart ReadyMix¹⁰
¹¹reagents (KAPA Biosystems, Inc. Wilmington, MA) was used to PCR the 16S rRNA¹¹
¹²gene. The PCR product amplicon size was verified using agarose gel electrophore-¹²
¹³sis. Concentration measurements were made after the initial 16S rRNA PCR, the¹³
¹⁴indexing PCR, and normalization steps. DNA concentration was measured using¹⁴
¹⁵the QuantIT Picogreen dsDNA Kit (Cat # P7589, ThermoFisher Scientific) and¹⁵
¹⁶fluorescent measurements were made with a Synergy2 Multi-Detection MicroPlate¹⁶
¹⁷Reader (BioTek Instruments, Inc, Winooski, VT).¹⁷

¹⁸ Initial PCR products were purified using 0.8X AMPure XP beads (Beckman Coul-¹⁸
¹⁹ter Genomics, Danvers, MA) following the manufacturer's protocol. After purifica-¹⁹
²⁰tion, the 192 samples were indexed using the Illumina Nextera XT index kits A²⁰
²¹and D (Illumina Inc., San Diego CA) and then purified using 1.12X AMPure XP²¹
²²beads. Prior to pooling purified sample concentration was normalized using Sequa-²²
²³Prep Normalization Plate Kit (Catalog n. A10510-01, Invitrogen Corp., Carlsbad,²³
²⁴CA), according to the manufacturer's protocol. Pooled library concentration was²⁴
²⁵checked using the Qubit dsDNA HS Assay Kit (Part# Q32851, Lot# 1735902,²⁵
²⁶ThermoFisher, Waltham, MA USA). Due to the low pooled amplicon library DNA²⁶
²⁷concentration, a modified protocol for low concentration libraries was used. The²⁷
²⁸library was run on an Illumina MiSeq, and base calls were made using Illumina²⁸
²⁹Real Time Analysis Software version 1.18.54. The sequence data was deposited in²⁹
³⁰the NCBI SRA archive under Bioproject PRJNA480312. Individual SRA run acces-³⁰
³¹sion numbers and metadata in Supplemental Table. Sequencing data quality control³¹
³²metrics for the 384 fastq sequence files (192 samples with forward and reverse reads)³²
³³were computed using the Bioconductor Rqc package [32, 33].³³

³⁴ Sequence data were processed using four bioinformatic pipelines: a *de-novo* clus-³⁴
³⁵tering method - Mothur [5], an open-reference clustering method - QIIME [7],³⁵
³⁶and a sequence inference method - DADA2 [8], and unclustered sequences as a³⁶
³⁷control. The code used to run the bioinformatic pipelines is available at [https:](https://github.com/nate-d-olson/mgtst_pipelines)³⁷
³⁸[/github.com/nate-d-olson/mgtst_pipelines](https://github.com/nate-d-olson/mgtst_pipelines).³⁸

³⁹ The Mothur pipeline follows the developer's MiSeq SOP [5, 23]. The pipeline was³⁹
⁴⁰run using Mothur version 1.37 (<http://www.mothur.org/>). We sequenced a larger⁴⁰
⁴¹16S rRNA region, with smaller overlap between the forward and reverse reads,⁴¹
⁴²than the 16S rRNA region the SOP was designed. Pipeline parameters modified to⁴²
⁴³account for difference in overlap are noted for individual steps below. The Make-⁴³
⁴⁴file and scripts used to run the Mothur pipeline are available [https://github.](https://github.com/nate-d-olson/mgtst_pipelines/blob/master/code/mothur)⁴⁴
⁴⁵[com/nate-d-olson/mgtst_pipelines/blob/master/code/mothur](https://github.com/nate-d-olson/mgtst_pipelines/blob/master/code/mothur). The Mothur⁴⁵
⁴⁶pipeline includes an initial preprocessing step where the forward and reverse reads⁴⁶

are trimmed and filtered using base quality scores and were merged into single contigs for each read pair. The following parameters were used for the initial contig filtering, no ambiguous bases, max contig length of 500 bp, and max homopolymer length of 8 bases. For the initial read filtering and merging step, low-quality reads were identified and filtered from the dataset based on the presence of ambiguous bases, failure to align to the SILVA reference database (V119, <https://www.arb-silva.de/>) [34], and identification as chimeras. Prior to alignment, the SILVA reference multiple sequence alignment was trimmed to the V3-V5 region, positions 6,388 and 25,316. Chimera filtering was performed using UChime (version v4.2.40) without a reference database [25]. OTU clustering was performed using the OptiClust algorithm with a clustering threshold of 0.97 [4]. The RDP classifier implemented in Mothur was used for taxonomic classification against the Mothur provided version of the RDP v9 training set [35].

The QIIME open-reference clustering pipeline for paired-end Illumina data was performed according to the online tutorial (Illumina Overview Tutorial (an IPython Notebook): open reference OTU picking and core diversity analyses, <http://qiime.org/tutorials/>) using QIIME version 1.9.1 [7]. Briefly, the QIIME pipeline uses fastq-join (version 1.3.1) to merge paired-end reads [36] and the Usearch algorithm [37] with Greengenes database version 13.8 with a 97% similarity threshold [38] was used for open-reference clustering.

DADA2, an R native pipeline was also used to process the sequencing data [8]. The pipeline includes a sequence inference step and taxonomic classification using the DADA2 implementation of the RDP naïve Bayesian classifier [35] and the SILVA database V123 provided by the DADA2 developers [34, <https://benjjneb.github.io/dada2/training.html>].

The unclustered pipeline was based on the Mothur *de-novo* clustering pipeline, where the paired-end reads were merged, filtered, and then dereplicated. Reads were aligned to the reference Silva alignment (V119, <https://www.arb-silva.de/>), and reads failing alignment were excluded from the dataset. Taxonomic classification of the unclustered sequences was performed using the same RDP classifier implemented in Mothur used for the *de-novo* pipeline. To limit the size of the dataset the most abundant 40,000 OTUs (comparable to the Mothur dataset), across all samples, were used as the unclustered dataset.

Qualitative Assessment

Artifactual Feature Proportion Our qualitative assessment evaluated features only observed in unmixed samples (PRE or POST) or only in titrations. The former we will refer to as unmixed-specific features and the latter we will refer to as titration-specific features (Fig. 1B). *Unmixed-* and *titration-specific* features can arise from errors in the PCR/sequencing, feature inference processes, or due to differences in sampling depth. To provide context for the artifactual feature proportion results count table sparsity was used (Fig. 1C). Sparsity is defined as the proportion of 0 valued cells in a matrix.

Hypothesis tests were used to determine if random sampling alone, here sequencing depth, could account for *unmixed-* and *titration-specific* features. p-values were adjusted for multiple comparisons using the Benjamini & Hochberg method [39].

¹For *unmixed-specific* features, a binomial test was used to evaluate if true feature¹
²relative abundance is less than the expected relative abundance. The binomial test²
³was infeasible for *titration-specific* features. Because the count table abundance³
⁴values for these features was 0 in the unmixed samples, their estimated probabil-⁴
⁵ity of occurrence π_{min} is equal to 0, and thus, the binomial test fails. Therefore,⁵
⁶we formulated a Bayesian hypothesis test for *titration-specific* features detailed by⁶
⁷equation (2). This Bayesian approach was used to evaluate if the true feature pro-⁷
⁸portion is less than the minimum detected proportion. Note that when assuming⁸
⁹equal priors, $P(\pi < \pi_{min}) = P(\pi > \pi_{min})$, (2) reduces to (3). We define π as the⁹
¹⁰true feature proportion, π_{min} the minimum detected proportion, C the expected¹⁰
¹¹feature counts, and C_{obs} the observed feature counts. Count values for C were sim-¹¹
¹²ulated using a beta prior (with varying alpha and beta values) for $\pi > \pi_{min}$ and¹²
¹³a uniform distribution for $\pi < \pi_{min}$. Higher values of alpha and beta will skew¹³
¹⁴the prior right and left respectively. Our Bayesian hypothesis tests (Eq. (3)) results¹⁴
¹⁵were largely unaffected by beta distribution parameterization (Fig. S4). π_{min} was¹⁵
¹⁶calculated using the mixture equation (1) where $q_{pre,j}$ and $q_{post,j}$ are $\min(\mathbf{Q}_{pre})$ ¹⁶
¹⁷and $\min(\mathbf{Q}_{post})$ across all features for a subject and pipeline. Our assumption is¹⁷
¹⁸that π is less than π_{min} for features not observed in unmixed samples. Artifacts not¹⁸
¹⁹explained by sequencing alone are likely errors in the sequence measurement and¹⁹
²⁰inference processes, and thus, false positives or negatives. ²⁰

$$\begin{aligned}
 & p = P(\pi < \pi_{min} | C \geq C_{obs}) \\
 & = \frac{P(C \geq C_{obs} | \pi < \pi_{min})P(\pi < \pi_{min})}{P(C \geq C_{obs} | \pi < \pi_{exp})P(\pi < \pi_{min}) + P(C \geq C_{obs} | \pi \geq \pi_{min})P(\pi \geq \pi_{min})} \quad (2) \\
 & p = \frac{P(C \geq C_{obs} | \pi < \pi_{min})}{P(C \geq C_{obs})} \quad (3)
 \end{aligned}$$

³⁰Quantitative Assessment

³¹For quantitative assessment, we compared observed relative abundance and log³¹
³²fold-changes to expected values derived from the titration experimental design.³²
³³Feature average relative abundance across PCR replicates was calculated using a³³
³⁴negative binomial model, and used as observed relative abundance values (*obs*) for³⁴
³⁵the relative abundance assessment. Average relative abundance values were used³⁵
³⁶to reduce PCR replicate outliers from biasing the assessment results. Equation (1)³⁶
³⁷and inferred θ values were used to calculate the expected relative abundance values³⁷
³⁸(*exp*). Relative abundance error rate is defined as $|exp - obs|/exp$. We developed³⁸
³⁹bias and variance metrics to assess feature performance. The feature-level bias and³⁹
⁴⁰variance metrics were defined as the median error rate and robust coefficient of⁴⁰
⁴¹variation ($RCOV = IQR/median$) respectively. ⁴¹

⁴²Log fold-change between samples in the titration series including PRE and POST⁴²
⁴³were compared to the expected log fold-change values to assess differential abun-⁴³
⁴⁴dance log fold-change estimates. Log fold-change estimates were calculated using⁴⁴
⁴⁵EdgeR [40, 41]. Expected log fold-change for feature j between titrations l and m ⁴⁵
⁴⁶is calculated using equation (4), where θ is the proportion of POST bacterial DNA⁴⁶

¹in a titration, and q is feature relative abundance. For features only present in PRE¹
²samples, the expected log fold-change is independent of the observed counts for the²
³unmixed samples and is calculated using (5). Features only observed in POST sam-³
⁴ples, *POST-specific*, expected log fold-change values can be calculated in a similar⁴
⁵manner. However, *POST-specific* features were rarely observed in more than one⁵
⁶titration and therefore were not suitable for use in our assessment. Due to a limited⁶
⁷number of *PRE-specific* features, both *PRE-specific* and *PRE-dominant* features⁷
⁸were used in the differential abundance assessment. *PRE-specific* features were de-⁸
⁹finied as features observed in all four PRE PCR replicates and not observed in any⁹
¹⁰of the POST PCR replicates and *PRE-dominant* features were also observed in all¹⁰
¹¹four PRE PCR replicates and observed in one or more of the POST PCR replicates¹¹
¹²with a log fold-change between PRE and POST samples greater than 5. ¹²

$$\log FC_{lm,j} = \log_2 \left(\frac{\theta_l q_{post,j} + (1 - \theta_l) q_{pre,i}}{\theta_m q_{post,j} + (1 - \theta_m) q_{pre,j}} \right) \quad (4)$$

$$\log FC_{lm,i} = \log_2 \left(\frac{1 - \theta_l}{1 - \theta_m} \right) \quad (5)$$

²¹Count Table Assessment Demonstration ²¹

²²Demonstrate framework by comparing the qualitative and quantitative assessment ²²
²³results across the three pipelines. We first characterized overall differences in the ²³
²⁴count tables produced by the three pipelines. This characterization included cal- ²⁴
²⁵culating the number of features, total abundance by sample, dropout-rate, and ²⁵
²⁶taxonomic composition. ²⁶

²⁸Qualitative Assessment ²⁸

²⁹For the qualitative assessment we compare the proportion of artifactual features. ²⁹
³⁰The artifactual feature proportion was defined as the proportion of *unmixed-* and ³⁰
³¹*titration-specific* features with abundance values that could not be explained by sam- ³¹
³²pling alone. These are PCR replicates with p-values less than 0.05 after multiple ³²
³³hypothesis test correction for the binomial and bayesian hypothesis tests described ³³
³⁴in the assessment framework methods section. We additionally used the count ta- ³⁴
³⁵ble sparsity values to draw conclusions regarding the mechanism responsible for ³⁵
³⁶different artifactual feature proportions. ³⁶

³⁸Quantitative Assessment ³⁸

³⁹Mixed-effects models were used to compare feature-level error rate bias and variance ³⁹
⁴⁰metrics across pipelines with subject as a random effect. Extreme feature-level error ⁴⁰
⁴¹rate bias and variance metric outliers were excluded from this analysis to minimize ⁴¹
⁴²biases due to poor model fit. Features with large bias and variance metrics, $1.5 \times IQR$ ⁴²
⁴³from the median, were deemed outliers. These outlier features were characterized ⁴³
⁴⁴independently in a separate analysis. ⁴⁴

⁴⁵We fit the following mixed effect model to test for differences in measurement bias ⁴⁵
⁴⁶across pipelines ⁴⁶

$$e_{ijk} = b + b_i + z_j + \epsilon_{ijk}$$

where e_{ijk} is the observed error across features and tritations k for pipeline i on individual j . b_i is a fixed term modeling the pipeline effect, z_j is a random effect (normally distributed with mean 0) capturing overall bias differences across individuals. We fit a similar model for differences in error variance across pipelines. We used estimated terms \hat{b}_i from the mixed effects model to test for pair-wise differences across pipelines. These multiple comparisons were performed with Tukey's HSD test. A one-sided alternative hypothesis was used to determine which pipelines had smaller feature-level error rate.

Declarations

Ethics approval and consent to participate
Not applicable.

Consent for publication
Not applicable.

Availability of data and material

Sequence data was deposited in the NCBI SRA archive under Bioproject PRJNA480312. Individual SRA run accession numbers and metadata in Supplemental Table. The code used to run the bioinformatic pipelines is available at https://github.com/nate-d-olson/mgtst_pipelines. Scripts used to analyze the data are available at https://github.com/nate-d-olson/mgtst_pub.

Competing interests

The authors declare that they have no competing interests.

Funding

This work was partially supported by National Institutes of Health (NIH) [NIH R01HG005220 to H.C.B.]

Authors' contributions

NDO, HCB, OCS, MS, and WT designed the experiment, SL and SH performed the laboratory work. NDO, HCB, MS, and DJB analyzed the data. NDO, DJB, and HCB wrote the manuscript. All authors provided feedback on manuscript drafts and approved the final manuscript.

Acknowledgements

The authors would like to thank the two anonymous reviewers, Mihai Pop, Scott Pine, Scott Jackson, Justin Zook, Nathan Swenson, and Prachi Kulkarni for feedback on manuscript drafts. Joseph Paulson and Justin Wagner provided helpful insight during the development of the project. Opinions expressed in this paper are the authors and do not necessarily reflect the policies and views of NIST, or affiliated venues. Certain commercial equipment, instruments, or materials are identified in this paper in order to specify the experimental procedure adequately. Such identification is not intended to imply recommendations or endorsement by NIST, nor is it intended to imply that the materials or equipment identified are necessarily the best available for the purpose. Official contribution of NIST; not subject to copyrights in USA.

Author details

¹Biosystems and Biomaterials Division, National Institute of Standards and Technology, 100 Bureau Dr., Gaithersburg, Maryland, 20899 USA. ²Center for Bioinformatics and Computational Biology, University of Maryland, College Park, 8314 Paint Branch Dr. College Park, Maryland, 20742 USA. ³University of Maryland Institute of Advanced Computer Studies, University of Maryland, College Park, 8223 Paint Branch Dr. College Park, Maryland, 20742 USA. ⁴Department of Epidemiology and Public Health, University of Maryland School of Medicine, 655 W. Baltimore Street, Baltimore, Maryland, 21201 USA. ⁵Department of Biomedical Engineering, Johns Hopkins University, 720 Rutland Ave., Baltimore, Maryland, 21205 USA. ⁶Joint Initiative for Metrology in Biology, 443 Via Ortega, Stanford, CA, 94305 USA. ⁷Department of Computer Science, University of Maryland, College Park, 8223 Paint Branch Dr. College Park, Maryland, 20742 USA.

References

1. Goodrich, J.K., Di Rienzi, S.C., Poole, A.C., Koren, O., Walters, W.A., Caporaso, J.G., Knight, R., Ley, R.E.: Conducting a microbiome study. *Cell* **158**(2), 250–262 (2014)
2. D'Amore, R., Ijaz, U.Z., Schirmer, M., Kenny, J.G., Gregory, R., Darby, A.C., Quince, C., Hall, N.: A comprehensive benchmarking study of protocols and sequencing platforms for 16S rRNA community profiling. *BMC Genomics* **17**, 1–40 (2016). doi:[10.1186/s12864-015-2194-9](https://doi.org/10.1186/s12864-015-2194-9)
3. Brooks, J.P., Edwards, D.J., Harwich, M.D., Rivera, M.C., Fettweis, J.M., Serrano, M.G., Reris, R.A., Sheth, N.U., Huang, B., Girerd, P., *et al.*: The truth about metagenomics: quantifying and counteracting bias in 16s rRNA studies. *BMC microbiology* **15**(1), 66 (2015)
4. Westcott, S.L., Schloss, P.D.: Opticlust, an improved method for assigning amplicon-based sequence data to operational taxonomic units. *mSphere* **2**(2) (2017)
5. Schloss, P.D., Westcott, S.L., Ryabin, T., Hall, J.R., Hartmann, M., Hollister, E.B., Lesniewski, R.A., Oakley, B.B., Parks, D.H., Robinson, C.J., *et al.*: Introducing mothur: open-source, platform-independent, community-supported software for describing and comparing microbial communities. *Applied and environmental microbiology* **75**(23), 7537–7541 (2009)
6. Rideout, J.R., He, Y., Navas-Molina, J.A., Walters, W.A., Ursell, L.K., Gibbons, S.M., Chase, J., McDonald, D., Gonzalez, A., Robbins-Pianka, A., Clemente, J.C., Gilbert, J.A., Huse, S.M., Zhou, H.-W., Knight, R., Caporaso, J.G.: Subsampled open-reference clustering creates consistent, comprehensive OTU definitions and scales to billions of sequences. *PeerJ* **2**, 545 (2014)
7. Caporaso, J.G., Kuczynski, J., Stombaugh, J., Bittinger, K., Bushman, F.D., Costello, E.K., Fierer, N., Peña, A.G., Goodrich, J.K., Gordon, J.I., Huttley, G.A., Kelley, S.T., Knights, D., Koenig, J.E., Ley, R.E., Lozupone, C.A., McDonald, D., Muegge, B.D., Pirrung, M., Reeder, J., Sevinsky, J.R., Turnbaugh, P.J., Walters, W.A., Widmann, J., Yatsunenko, T., Zaneveld, J., Knight, R.: Qiime allows analysis of high-throughput community sequencing data. *Nature Methods* **7**, 335 (2010). Correspondence
8. Callahan, B.J., McMurdie, P.J., Rosen, M.J., Han, A.W., Johnson, A.J.A., Holmes, S.P.: Dada2: High-resolution sample inference from illumina amplicon data. *Nature Methods* **13**, 581–583 (2016). doi:[10.1038/nmeth.3869](https://doi.org/10.1038/nmeth.3869)
9. Bokulich, N.A., Rideout, J.R., Mercurio, W.G., Shiffer, A., Wolfe, B., Maurice, C.F., Dutton, R.J., Turnbaugh, P.J., Knight, R., Caporaso, J.G.: mockrobiota: a public resource for microbiome bioinformatics benchmarking. *mSystems* **1**(5), 00062–16 (2016)
10. Kopylova, E., Navas-molina, J.A., Mercier, C., Xu, Z.: Open-Source Sequence Clustering Methods Improve the State Of the Art. *mSystems* **1**(1), 1–16 (2014). doi:[10.1128/mSystems.00003-15](https://doi.org/10.1128/mSystems.00003-15). Editor
11. Huse, S.M., Welch, D.M., Morrison, H.G., Sogin, M.L.: Ironing out the wrinkles in the rare biosphere through improved OTU clustering. *Environmental microbiology* **12**(7), 1889–98 (2010). doi:[10.1111/j.1462-2920.2010.02193.x](https://doi.org/10.1111/j.1462-2920.2010.02193.x)
12. Costea, P.I., Zeller, G., Sunagawa, S., Pelletier, E., Alberti, A., Levenez, F., Tramontano, M., Driessen, M., Herczeg, R., Jung, F.-E., Kultima, J.R., Hayward, M.R., Coelho, L.P., Allen-Vercoe, E., Bertrand, L., Blaut, M., Brown, J.R.M., Carton, T., Cools-Portier, S., Daigneault, M., Derrien, M., Druesne, A., de Vos, W.M., Finlay, B.B., Flint, H.J., Guarner, F., Hattori, M., Heilig, H., Luna, R.A., van Hylckama Vlieg, J., Junick, J., Klymiuk, I., Langella, P., Le Chatelier, E., Mai, V., Manichanh, C., Martin, J.C., Mery, C., Morita, H., O'Toole, P.W., Orvain, C., Patil, K.R., Penders, J., Persson, S., Pons, N., Popova, M., Salonen, A., Saulnier, D., Scott, K.P., Singh, B., Slezak, K., Veiga, P., Versalovic, J., Zhao, L., Zoetendal, E.G., Ehrlich, S.D., Dore, J., Bork, P.: Towards standards for human fecal sample processing in metagenomic studies. *Nat. Biotechnol.* **35**, 1069 (2017)
13. Olson, N.D., Morrow, J.B.: DNA extract characterization process for microbial detection methods development and validation. *BMC Res. Notes* **5**, 668 (2012)
14. Klindworth, A., Pruesse, E., Schweer, T., Peplies, J., Quast, C., Horn, M., Glöckner, F.O.: Evaluation of general 16s ribosomal rna gene pcr primers for classical and next-generation sequencing-based diversity studies. *Nucleic acids research*, 808 (2012)
15. Gohl, D.M., Vangay, P., Garbe, J., MacLean, A., Hauge, A., Becker, A., Gould, T.J., Clayton, J.B., Johnson, T.J., Hunter, R., Knights, D., Beckman, K.B.: Systematic improvement of amplicon marker gene methods for increased accuracy in microbiome studies. *Nat. Biotechnol.* (2016)
16. Pinto, A.J., Raskin, L.: PCR biases distort bacterial and archaeal community structure in pyrosequencing datasets. *PLoS One* **7**(8), 43093 (2012)
17. Hansen, M.C., Tolker-Nielsen, T., Givskov, M., Molin, S.: Biased 16S rDNA PCR amplification caused by interference from DNA flanking the template region. *FEMS Microbiol. Ecol.* **26**(2), 141–149 (1998)
18. McMurdie, P.J., Holmes, S.: Waste not, want not: why rarefying microbiome data is inadmissible. *PLoS Comput. Biol.* **10**(4), 1003531 (2014)
19. Parsons, J., Munro, S., Pine, P.S., McDaniel, J., Mehaffey, M., Salit, M.: Using mixtures of biological samples as process controls for rna-sequencing experiments. *BMC genomics* **16**(1), 708 (2015)
20. Pine, P.S., Rosenzweig, B.A., Thompson, K.L.: An adaptable method using human mixed tissue ratiometric controls for benchmarking performance on gene expression microarrays in clinical laboratories. *BMC biotechnology* **11**(1), 38 (2011)
21. Thompson, K.L., Rosenzweig, B.A., Pine, P.S., Retief, J., Turpaz, Y., Afshari, C.A., Hamadeh, H.K., Damore, M.A., Boedigheimer, M., Blomme, E., *et al.*: Use of a mixed tissue rna design for performance assessments on multiple microarray formats. *Nucleic acids research* **33**(22), 187–187 (2005)
22. Pop, M., Paulson, J.N., Chakraborty, S., Astrovskaia, I., Lindsay, B.R., Li, S., Bravo, H.C., Harro, C., Parkhill, J., Walker, A.W., *et al.*: Individual-specific changes in the human gut microbiota after challenge with enterotoxigenic *escherichia coli* and subsequent ciprofloxacin treatment. *BMC genomics* **17**(1), 1 (2016)
23. Kozich, J.J., Westcott, S.L., Baxter, N.T., Highlander, S.K., Schloss, P.D.: Development of a dual-index sequencing strategy and curation pipeline for analyzing amplicon sequence data on the miseq illumina

- sequencing platform. *Applied and environmental microbiology* **79**(17), 5112–5120 (2013)
224. Walters, W., Hyde, E.R., Berg-Lyons, D., Ackermann, G., Humphrey, G., Parada, A., Gilbert, J.A., Jansson, J.K., Caporaso, J.G., Fuhrman, J.A., Apprill, A., Knight, R.: Improved bacterial 16S rRNA gene (v4 and v4-5) and fungal internal transcribed spacer marker gene primers for microbial community surveys. *mSystems* **1**(1) (2016)
25. Edgar, R.C., Haas, B.J., Clemente, J.C., Quince, C., Knight, R.: Uchime improves sensitivity and speed of chimera detection. *Bioinformatics* **27**(16), 2194–2200 (2011)
26. Edgar, R.C.: UNCross2: identification of cross-talk in 16S rRNA OTU tables (2018)
27. D'Amore, R., Ijaz, U.Z., Schirmer, M., Kenny, J.G., Gregory, R., Darby, A.C., Shakya, M., Podar, M., Quince, C., Hall, N.: A comprehensive benchmarking study of protocols and sequencing platforms for 16S rRNA community profiling. *BMC Genomics* **17**, 55 (2016)
28. Sze, M.A., Schloss, P.D.: The impact of dna polymerase and number of rounds of amplification in pcr on 16s rRNA gene sequence data. *bioRxiv* (2019). doi:[10.1101/565598](https://doi.org/10.1101/565598). <https://www.biorxiv.org/content/early/2019/03/04/565598.full.pdf>
29. Wright, E.S., Yilmaz, L.S., Ram, S., Gasser, J.M., Harrington, G.W., Noguera, D.R.: Exploiting extension bias in polymerase chain reaction to improve primer specificity in ensembles of nearly identical dna templates. *Environmental microbiology* **16**(5), 1354–1365 (2014)
30. Harro, C., Chakraborty, S., Feller, A., DeNearing, B., Cage, A., Ram, M., Lundgren, A., Svennerholm, A.-M., Bourgeois, A.L., Walker, R.I., *et al.*: Refinement of a human challenge model for evaluation of enterotoxigenic *Escherichia coli* vaccines. *Clinical and Vaccine Immunology* **18**(10), 1719–1727 (2011)
31. Yang, B., Wang, Y., Qian, P.-Y.: Sensitivity and correlation of hypervariable regions in 16s rRNA genes in phylogenetic analysis. *BMC bioinformatics* **17**(1), 1 (2016)
32. Souza, W., Carvalho, B.: Rqc: Quality Control Tool for High-Throughput Sequencing Data. (2017). R package version 1.10.2. <https://github.com/labcb/Rqc>
33. Huber, W., Carey, J., V., Gentleman, R., Anders, S., Carlson, M., Carvalho, S., B., Bravo, C., H., Davis, S., Gatto, L., Girke, T., Gottardo, R., Hahne, F., Hansen, D., K., Irizarry, A., R., Lawrence, M., Love, I., M., MacDonald, J., Obenchain, V., Ole's, K., A., Pag'es, H., Reyes, A., Shannon, P., Smyth, K., G., Tenenbaum, D., Waldron, L., Morgan, M.: Orchestrating high-throughput genomic analysis with Bioconductor. *Nature Methods* **12**(2), 115–121 (2015)
34. Quast, C., Pruesse, E., Yilmaz, P., Gerken, J., Schweer, T., Yarza, P., Peplies, J., Glöckner, F.O.: The SILVA ribosomal RNA gene database project: improved data processing and web-based tools. *Nucleic acids research* **41**(D1), 590–596 (2012)
35. Wang, Q., Garrity, G.M., Tiedje, J.M., Cole, J.R.: Naive bayesian classifier for rapid assignment of rRNA sequences into the new bacterial taxonomy. *Applied and environmental microbiology* **73**(16), 5261–5267 (2007)
36. Aronesty, E.: ea-utils: Command-line tools for processing biological sequencing data. *Expression Analysis*, Durham, NC (2011)
37. Edgar, R.C.: Search and clustering orders of magnitude faster than blast. *Bioinformatics* **26**(19), 2460–2461 (2010)
38. DeSantis, T.Z., Hugenholtz, P., Larsen, N., Rojas, M., Brodie, E.L., Keller, K., Huber, T., Dalevi, D., Hu, P., Andersen, G.L.: Greengenes, a chimera-checked 16s rRNA gene database and workbench compatible with arb. *Applied and environmental microbiology* **72**(7), 5069–5072 (2006)
39. Benjamini, Y., Hochberg, Y.: Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the royal statistical society. Series B (Methodological)*, 289–300 (1995)
40. Robinson, M.D., McCarthy, D.J., Smyth, G.K.: edgeR: a bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* **26**(1), 139–140 (2010)
41. McCarthy, D.J., Chen, Y., Smyth, G.K.: Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Res.* **40**(10), 4288–4297 (2012)