

Taxonomic classification method for metagenomics based on core protein families with Core-Kaiju

Anna Tovo^{1,2}, Peter Menzel³, Anders Krogh⁴, Marco Cosentino Lagomarsino^{5,6} and Samir Suweis^{1,7,*}

¹Physics and Astronomy Department, LIPh Lab, University of Padova, Via Marzolo 8, 35131 Padova, Italy, ²Mathematics Department, University of Padova, via Trieste 63, 35121 Padova, Italy, ³Labor Berlin Charité Vivantes GmbH, Sylter Str. 2, 13353 Berlin, Germany, ⁴Department of Computer Science, University of Copenhagen, Universitetsparken 1, DK-2100 Copenhagen, Denmark, ⁵IFOM, FIRG Institute of Molecular Oncology, Via Adamello 16, 20143 Milan, Italy, ⁶Physics Department, University of Milan, and I.N.F.N., Via Celoria 16, 20133 Milan, Italy and ⁷Padova Neuroscience Center, University of Padova, Via Orus 2/B, 35131 Padova, Italy

Received January 14, 2020; Revised June 12, 2020; Editorial Decision June 22, 2020; Accepted June 24, 2020

ABSTRACT

Characterizing species diversity and composition of bacteria hosted by biota is revolutionizing our understanding of the role of symbiotic interactions in ecosystems. Determining microbiomes diversity implies the assignment of individual reads to taxa by comparison to reference databases. Although computational methods aimed at identifying the microbe(s) taxa are available, it is well known that inferences using different methods can vary widely depending on various biases. In this study, we first apply and compare different bioinformatics methods based on 16S ribosomal RNA gene and shotgun sequencing to three mock communities of bacteria, of which the compositions are known. We show that none of these methods can infer both the true number of taxa and their abundances. We thus propose a novel approach, named Core-Kaiju, which combines the power of shotgun metagenomics data with a more focused marker gene classification method similar to 16S, but based on emergent statistics of core protein domain families. We thus test the proposed method on various mock communities and we show that Core-Kaiju reliably predicts both number of taxa and abundances. Finally, we apply our method on human gut samples, showing how Core-Kaiju may give more accurate ecological characterization and a fresh view on real microbiomes.

INTRODUCTION

Modern high-throughput genome sequencing techniques revolutionized ecological studies of microbial communities

at an unprecedented range of taxa and scales (1–5). It is now possible to massively sequence genomic DNA directly from incredibly diverse environmental samples (3,6) and gain novel insights about structure and metabolic functions of microbial communities.

One major biological question is the inference of the composition of a microbial community, that is, the relative abundances of the sampled organisms. In particular, the impact of microbial diversity and composition for the maintenance of human health is increasingly recognized (7–10). Indeed, several studies suggest that the disruption of the normal microbial community structure, known as dysbiosis, is associated with diseases ranging from localized gastroenterologic disorders (11) to neurologic illnesses (12). However, it is impossible to define dysbiosis without first establishing what ‘normal microbial community structure’ means within the healthy human microbiome. To this purpose, the Human Microbiome Project has analysed the largest cohort and set of distinct, clinically relevant body habitats (13), characterizing the ecology of healthy human-associated microbial communities. However, there are several critical aspects. The study of the structure, function and diversity of the human microbiome has revealed that even healthy individuals differ remarkably in the contained species and their abundances. Much of this diversity remains unexplained, although diet, environment, host genetics and early microbial exposure have all been implicated. Characterizing a microbial community implies the classification of species/genera composition within the sampled community, which in turn requires the assignment of sequencing reads to taxa, usually by comparison to a reference database. Although computational methods aimed at identifying the microbe(s) taxa have an increasingly long history within bioinformatics (14–16), it is well known that inference based on 16S ribosomal RNA (rRNA) or shotgun sequencing vary widely (17). Moreover, even if data are obtained via the same experimen-

*To whom correspondence should be addressed. Tel: +39 049 827 7174; Email: samir.suweis@unipd.it

tal protocol, the usage of different computational methods or algorithm variants may lead to different results in the taxonomic classification. The two main experimental approaches for analyzing the microbiomes are based on 16S rRNA gene amplicon sequencing and whole genome shotgun sequencing (metagenomics).

Sequencing of amplicons from a region of the 16S rRNA gene is a common approach used to characterize microbiomes (18,19) and many analysis tools are available (see Materials and Methods section). Besides the biases in the experimental protocol, a major issue with 16S amplicon-sequencing is the variance of copy numbers of the 16S genes between different taxa. Therefore, abundances inferred by read counts of the amplicons should be properly corrected by taking into account the copy number of the different genera detected in the sample (3,20,21). However, the average number of 16S rRNA copies is only known for a restricted selection of bacterial taxa. As a consequence, different algorithms have been proposed to infer from data the copy number of those taxa for which this information is not available (18,22).

In contrast, whole genome shotgun sequencing of all the DNA present in a sample can inform about both diversity and abundance as well as metabolic functions of the species in the community (23). The accuracy of shotgun metagenomics species classification methods varies widely (24). In particular, these methods can typically result in a large number of false positive predictions, depending on the used sequence comparison algorithm and its parameters. For example in *k*-mer based methods as Kraken (25) and Kraken2 (26) the choice of *k* determines sensitivity and precision of the classification, such that sensitivity increases and precision decreases with increasing values for *k*, and vice versa. As we will show, false positive predictions often need to be corrected heuristically by removing all taxa with abundance below a given arbitrary threshold (see Materials and Methods section for an overview on different algorithms of taxonomy classification).

We highlight that the protocols for 16S-amplicons and shotgun methods are different and each has their own batch effects. Importantly, while shotgun taxonomic analysis gives classification results at species-level, 16S taxonomic profilers most often need to stop at the genus level. However, in the end, both aim at answering to the same question: ‘what are the relative abundances of taxa in the sample?’ Therefore, it makes sense methodologically to compare their answers against the same community. To do that, it is possible to aggregate lower level (e.g. species) counts towards higher levels (e.g. genus), as it has been done in many benchmarks studies before (see, e.g. (17,25,27,28)). In fact, several studies have performed comparisons of taxa inferred from 16S amplicon and shotgun sequencing data, with samples ranging from humans to studies of water and soil. Logares and collaborators (29) studied communities of bacteria marine plankton and found that shotgun approaches had an advantage over amplicons, as they rendered more truthful community richness and evenness estimates by avoiding PCR biases, and provided additional functional information. Chan et al. (30) analyzed thermophilic bacteria in hot spring water and found that amplicon and shotgun sequencing al-

lowed for comparable phylum detection, but shotgun sequencing failed to detect three phyla. In another study (31) 16S rRNA and shotgun methods were compared in classifying community bacteria sampled from freshwater. Taxonomic composition of each 16S rRNA gene library was generally similar to its corresponding metagenome at the phylum level. At the genus level, however, there was a large amount of variation between the 16S rRNA sequences and the metagenomic contigs, which had a ten-fold resolution and sensitivity for genus diversity. More recently Jovel *et al.* (27) compared bacteria communities from different microbiomes (human, mice) and also from mock communities. They found that shotgun metagenomics offered a greater potential for identification of strains, which however still remained unsatisfactory. It also allowed increased taxonomic and functional resolution, as well as the discovery of new genomes and genes.

While shotgun metagenomics has certain advantages over amplicon-sequencing, its higher price point is still prohibitive for many applications. Therefore, amplicon sequencing remains the go-to established cost-effective tool to the taxonomic composition of microbial communities. In fact, the usage of the 16S rRNA-gene as a universal marker throughout the entire bacterial kingdom made it easy to collect sequence information from a wide distribution of taxa, which is yet unmatched by whole genome databases. Several curated databases exist to date, with SILVA (32,33), GreenGenes (34,35) and Ribosomal Database Project (RDP) (36) being the most prominent. Additionally, NCBI also provides a curated collection of 16S reference sequences in its Targeted Loci project (<https://www.ncbi.nlm.nih.gov/refseq/targetedloci/>).

When benchmarking protocols for taxonomic classification from real samples of complex microbiomes, the ‘ground truth’ of the contained taxa and their relative abundances is not known (see (27)). Therefore, the use of mock communities or simulated datasets remains as basis for a robust comparative evaluation of a method prediction accuracy. In the first part of this work, we apply three widely used taxonomic classifiers for metagenomics, Kaiju (28), Kraken2 (26) and MetaPhlan2 (37), and two common methods for analyzing 16S-amplicon sequencing data, DADA2 (38) and QIIME2 (39) to three small mock communities of bacteria, of which we know the exact composition (27). We show that 16S rRNA data efficiently allow to detect the number of taxa, but not their abundances, while shotgun metagenomics as Kaiju and Kraken2 give a reliable estimate of the most abundant genera, but the nature of the algorithms makes them predict a very large number of false-positive taxa.

The central contribution of this work is thus to develop a method to overcome the above limitations. In particular, we propose an updated version of Kaiju, which combines the power of shotgun metagenomics data with a more focused marker gene classification method, similar to 16S rRNA, but based on core protein domain families (40–43) from the PFAM database (44).

Our criterion for choosing the set of marker domain families is that we uncover the existence of a set of core families that are typically at most present in one or very few copies

per genome, but together cover uniquely all 8116 bacteria species in the PFAM database with an overall quite short sequence. Using presence of these core PFAMs (mostly related to ribosomal proteins) as a filter criterion allows for detecting the correct number of taxa in the sample. We tested our approach in a protocol called ‘Core-Kaiju’ and show that it has a higher accuracy than other classification methods not only on the three small mock communities, but also on intermediate and highly biodiverse mock communities designed for the 1st Critical Assessment of Metagenome Interpretation (CAMI) challenge (45). In fact, we will show how in all these cases Core-Kaiju overcomes, for the most part, the problem of false-positive genera and accurately predicts the abundances of the different detected taxa. We finally apply our novel pipeline to classify microbial genera in the human gut from the Human Microbiome Project (HMP) (see <https://www.hmpdacc.org/hmp/HM16STR/>) dataset, showing how Core-Kaiju may allow for a more accurate biodiversity characterization of real microbial communities, thus putting the basis for more solid analysis in microbiomes.

MATERIALS AND METHODS

Taxonomic classification: amplicon versus whole genome sequencing

Many computational tools are available for the analysis of both amplicon and shotgun sequencing data (25,26,28,37–39,46).

One of the differences among the several software for 16S rRNA analysis, is on how the next-generation sequencing error rate per nucleotide is taken into account, when associating each sampled 16S sequence read to taxa. Indeed, errors along the nucleotide sequence could lead to an inaccurate taxon identification and, consequently, to misleading diversity statistics.

The traditional approach to overcome this problem is to cluster amplicon sequences into the so-called operational taxonomic units (OTUs), which are based on an arbitrary shared similarity threshold usually set up equal to 97% for classification at the genus level. Of course, in this way, these approaches lead to a reduction of the phylogenetic resolution, since gene sequences below the fixed threshold cannot be distinguished one from the other.

That is why, sometimes, it may be preferable to work with exact amplicon sequence variants (ASVs), i.e. sequences recovered from a high-throughput marker gene analysis after the removal of spurious sequences generated during PCR amplification and/or sequencing techniques. The next step in these approaches is to compare the filtered sequences with reference libraries as those cited above. In this work, we chose to conduct the analyses with the following two open-source platforms: DADA2 (38) and QIIME2 (39). DADA2 is an R-package optimized to process large datasets (from 10s of millions to billions of reads) of amplicon sequencing data with the aim of inferring the ASVs from one or more samples. Once the spurious 16S rRNA gene sequences have been recovered, DADA2 allowed for the comparison with both SILVA, GreenGenes and RDP libraries. We performed the analyses for all the three possible choices. QIIME2 is

another widely used bioinformatic platform for the exploration and analysis of microbial data which allows, for the sequence quality control step, to choose between different methods. For our comparisons, we performed this step by using Deblur (47), a novel sub-operational-taxonomic-unit approach which exploits information on error profiles to recover error-free 16S rRNA sequences from samples.

As shown in (27), where different amplicon sequencing methods are tested on both simulated and real data and the results are compared to those obtained with metagenomic pipelines, the whole genome approach resulted to outperform the previous ones in terms of both number of identified strains, taxonomic and functional resolution and reliability on estimates of microbial relative abundance distribution in samples.

Similar comparisons have also been performed with analogous results in (29,30,46,48) (see (17) for a comprehensive summary of studies comparing different sequencing approaches and bioinformatic platforms).

Standard widespread taxonomic classification algorithms for metagenomics (e.g. Kraken (25) and Kraken2 (26)) extract all contained k -mers (all the possible strings of length k that are contained in the whole metagenome) from the sequencing reads and compare them with index of a genome database. However, the choice of the length k highly influences the classification, since, when k is too large, it is easy not to find a correspondence in reference database, whereas if k is too small, reads may be wrongly classified. Recently, a novel approach has been proposed for the classification of shotgun data based on sequence comparison to a reference database comprising protein sequences, which are much more conserved with respect to nucleotide sequences (28). Kaiju indexes the reference database using the Borrows-Wheeler-Transform (BWT), and translated sequencing reads are searched in the BWT using maximum exact matches, optionally allowing for a certain number of mismatches via a greedy heuristic approach. It has been shown (28) that Kaiju is able to classify more reads in real metagenomes than nucleotide-based k -mers methods. Therefore, previous studies on the community composition and structure of microbial communities in the human can be actually very biased by previous metagenomic analysis that were missing up to 90% of the reconstructed species (i.e. most of the species they found were not present in the gene catalog). We therefore chose to work with Kaiju (with MEM option (28)) for our taxonomic analysis. Although it resulted to give better estimates of sample biodiversity composition with respect to amplicon sequencing techniques, we found that it generally overestimates the number of genera actually present in our community (see Results section) of two magnitude orders, i.e. there is a long tail of low abundant false-positive taxa. To overcome this, we implemented a new release of the program, Core-Kaiju, which contains an additional preliminary step where reads sequences are firstly mapped against a newly protein reference library we created containing the amino-acid sequence of proteomes’ core PFAMs (see following section). We also compared standard Kaiju and Core-Kaiju results with those obtained via Kraken2 and via another widely used program for shotgun data analysis, MetaPhlan2 (37,46).

Characterization of the core PFAM families

After downloading the PFAM database (version 32.0), we selected only bacterial proteomes and we tabulated the data into a $F \times P$ matrix, where each column represented a different proteome and each row a different protein domain. In particular, our database consisted of $P = 8116$ bacterial proteomes and $F = 11\,286$ protein families. In each matrix entry (f, p) , we inserted the number of times the f family recurred in proteins of the p proteome, $n_{f,p}$. By summing up over the p column, one can get the proteome length, i.e. the total number of families of which it is constituted, which we will denote with l_p . Similarly, if we sum up over the f row, we get the family abundance, i.e. the number of times the f family appears in the PFAM database, which we call a_f . Figure 1 shows the frequency histogram of the proteome sizes (left panel) and of the family abundances (right panel). Our primary goal was to find the so-called core families (49), i.e. the protein domains which are present in the overwhelming majority of the bacterium proteomes but occurring just few times in each of them (41,50). In order to analyze the occurrences of PFAM in proteomes, we converted the original $F \times P$ matrix into a binary one, giving information on whether each PFAM was present or not in each proteome. In the left panel of Figure 2, we inserted the histogram of the family occurrences, which displays the typical *u-shape*, already observed in literature (43,51–53): a huge number of families are present in only few proteomes (first peak in the histogram), whilst another smaller peak occurs at large values, meaning that there are also a percentage of domains occurring in almost all the proteomes. In the right panel, we show the plot of the number of rare PFAM (having abundance less or equal to four in each proteome) versus the percentage of proteomes in which they have been found. We thus selected the PFAMs found in more than 90% of the proteomes and such that $\max_p n_{f,p} = 4$ (see Zoom 2 panel of Figure 2).

Since we wish to have at least one representative core PFAM for each proteome in the database, we checked whether with these selected core families we could ‘cover’ all bacteria. Unfortunately, none of them resulted to be present in proteomes 479430 and 1609106, corresponding to *Actinospica robiniae* DSM 44927 and *Streptomyces* sp. NRRL B-1568, respectively. We therefore looked for the most prevalent PFAM(s) present in such proteomes. We found that PFAM PF08338, occurring in 43% of the proteomes, was present in both *Actinospica robiniae* and *Streptomyces* and we therefore add it to our core-PFAM list. Eventually, in order to minimize the number of PFAMs to work with (and related computational cost), we considered in our final core-PFAM list only the minimum number of domains through which we were able to cover the whole list of proteomes of the databases. In particular, the selected core protein domains for bacteria proteomes are the ten PFAMs PF00453, PF00572, PF01029, PF01649, PF01795, PF03947, PF08338, PF09285 and PF17136 (see Table 1).

Principal coordinate analysis. In order to explore whether the expression of the core PFAM protein domains are correlated with taxonomy, we did the following. First, we downloaded from the UniProt database (54) the amino acid se-

quence of each PFAM along the different proteomes (see Supplementary Figures S1 and S2, for details). Their averaged (over proteomes) sequence lengths L resulted to be highly picked around specific values ranging from $L = 46$ to $L = 297$ (see Supplementary Figure S3, for the corresponding frequency histograms).

Second, for each family we computed the Damerau–Levenshtein (DL) distance between all its corresponding DNA sequences. DL measures the edit distance between two strings in terms of the minimum number of allowed operations needed to modify one string to match the other. Such operations include insertions, deletions/substitutions of single characters and transposition of two adjacent characters, which are common errors occurring during DNA polymerase. This analogy makes the DL distance a suitable metric for the variation between protein sequences. By simplicity and to have a more immediate insight, we conducted the analysis only for sequence points corresponding to the five most abundant phyla, i.e. Proteobacteria, Firmicutes, Actinobacteria, Bacteroidetes and Cyanobacteria.

After computing the DL distance matrices between all the amino-acid sequences of each PFAMs along proteomes, we performed the Multi Dimensional Scaling (MDS) or Principal Coordinate Analysis (PCoA) on the DL distance matrix. This step allow us to reduce the dimensionality of the space describing the distances between all pairs of core PFAMs of the different taxa and visualize it in a two dimensional space. In the last two columns of Table 2 we inserted the percentage of the variance explained by the first two principal coordinates for the ten different core families, where the first one ranges from 3.3 to 12.1% and the second one from 2.4 to 7.7%. We then plotted the sequence points into the new principal coordinate space, coloring them by phyla. In general, we observed a two-case scenario. For some families as PF03883 (see Figure 3, left panel), Actinobacteria and Proteobacteria sequences are grouped in one or two highly visible clusters each, whereas the other three phyla do not form well distinguished structures, being their sequence points close one another, especially for Cyanobacteria and Firmicutes. For other families as PF01196 (see Figure 3, left panel), all five phyla result to be clustered, suggesting a higher correlation between taxonomy and amino-acid sequences (see Supplementary Figure S4, for the other core families graphics). These results suggest that some core families (e.g. ribosomal ones) are phyla dependent, while other are not directly correlated with taxa.

Mock bacteria communities

We started by testing shotgun versus 16S taxonomic pipelines on three small artificial bacterial communities generated by Jovel *et al.* (27), whose raw data are publicly available (Sequence Read Archive (SRA) portal of NCBI, accession number SRP059928). These mock populations contain DNA from eleven species belonging to seven genera: *Salmonella enterica*, *Streptococcus pyogenes*, *Escherichia coli*, *Lactobacillus helveticus*, *Lactobacillus delbrueckii*, *Lactobacillus plantarum*, *Clostridium sordelli*, *Bacteroides thetaiotaomicron*, *Bacteroides vulgatus*, *Bifidobacterium breve* and *Bifidobacterium animalis*. For the taxo-

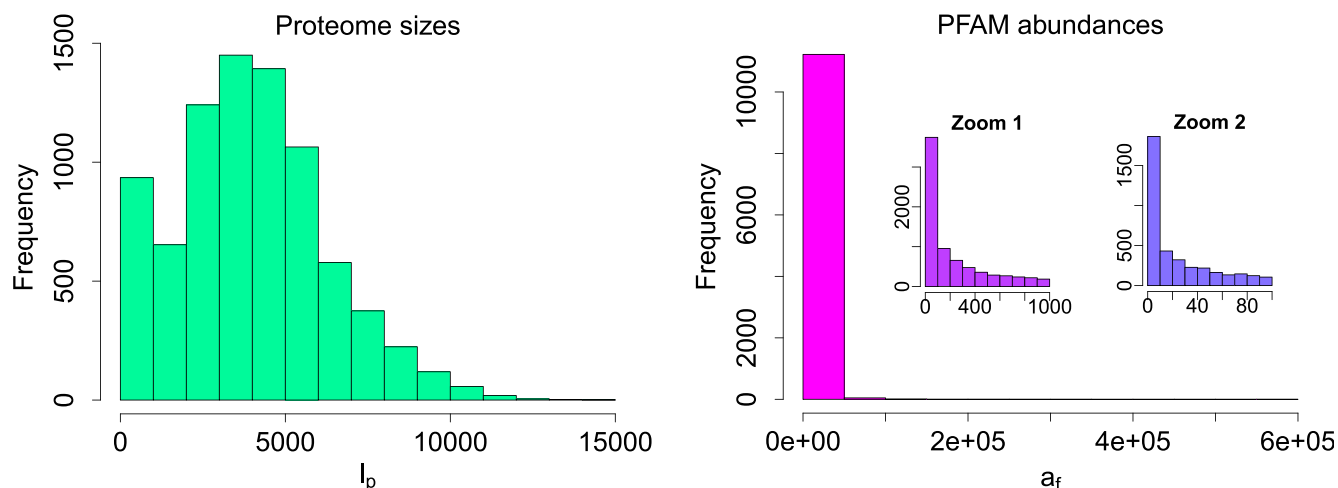


Figure 1. Proteome sizes and families abundances in PFAM database. On the left panel: frequency histogram of proteome lengths l_p (total number of families of which a proteome p is composed). On the right panel: frequency histogram of family abundances, a_f (number of times a PFAM f appears along a proteome).

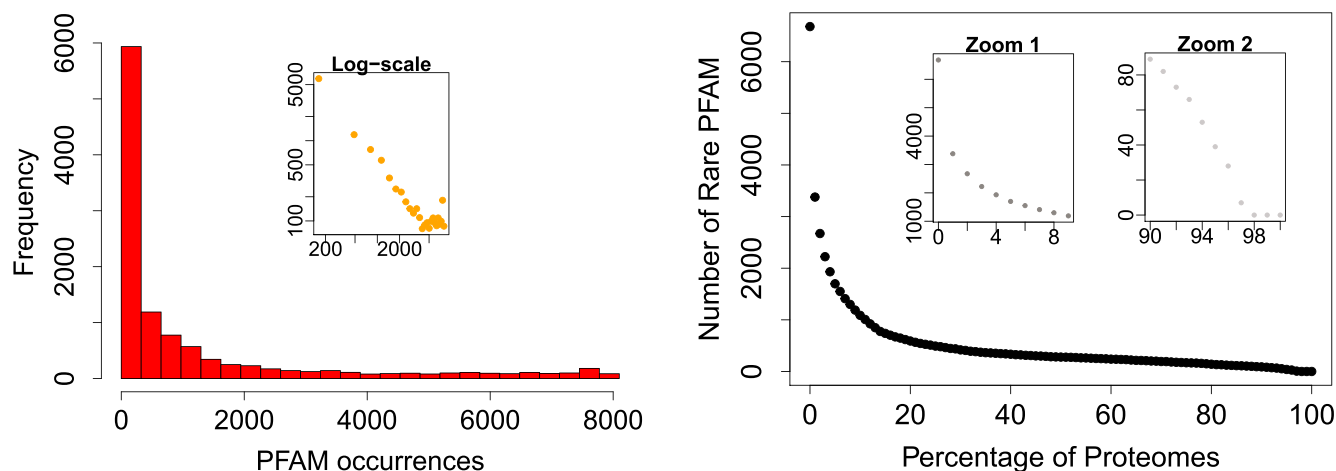


Figure 2. PFAM occurrences along proteomes. On the left panel: frequency histogram of family occurrences (number of proteomes in which a PFAM is contained). On the right panel: number of families with occurrence at most four versus the percentage of proteomes in which they are contained.

Table 1. Core PFAMs identity number and corresponding function in proteomes

PFAM ID	Function
PF00453	Ribosomal protein L20
PF00572	Ribosomal protein L13
PF01029	NusB family (involved in the regulation of rRNA biosynthesis by transcriptional antitermination)
PF01196	Ribosomal protein L17
PF01649	Ribosomal protein S20 (Bacterial ribosomal protein S20 interacts with 16S rRNA)
PF01795	MraW methylase family (SAM dependent methyltransferases)
PF03947	Ribosomal Proteins L2, C-terminal domain
PF08338	Domain of unknown function (DUF1731)
PF09285	EF-P (elongation factor P) translation factor required for efficient peptide bond synthesis on 70S ribosomes
PF17136	Ribosomal proteins 50S L24/mitochondrial 39S L24

nomic analysis at the genus level through 16S amplicon sequencing, we evaluated the performance of DADA2 (38) and QIIME2 pipelines (39). In particular, as shown in (27), QIIME2 produced more reliable results in terms of relative abundance of bacteria for all three mock communities when compared to Mothur (55), another widely used 16S

pipeline, and to the MiSeq Reporter v2.5, a software developed by Illumina to analyze MiSeq instrument output data.

As for shotgun libraries, we tested the standard Kaiju (28), Kraken2 (26), the improved version of Kraken (25), and MetaPhlAn2 (37), the improved version of MetaPhlAn (46). This latter relies on unique clade-specific marker genes

Table 2. Prevalence, maximal/total occurrences and principal coordinates of PFAM core families. We inserted, for each core family (PFAM ID, first column), the percentage of proteomes in which it appears (prevalence, second column), the maximum number of times it occurs in one proteome (maximal occurrence, third column), the total number of times it is found among proteomes in the PFAM database (total occurrence, fourth column) and the percentage of variance explained by the first two coordinates (PCo1 and PCo2, last two columns) when MDS is performed on sequences belonging to the five most abundant phyla (see Figure 3)

PFAM ID	Prevalence	Maximal occurrence	Total occurrence	PCo1	PCo2
PF00453	95%	3	7786	10.6%	6.6%
PF00572	97%	3	7897	5.4%	5.1%
PF01029	96%	4	12991	3.9%	2.4%
PF01196	97%	3	7888	12.1%	5.7%
PF01649	94%	3	7715	6.1%	4.6%
PF01795	96%	4	8113	5.2%	4.9%
PF03947	97%	4	7886	8.2%	7.7%
PF08338	43%	4	4267	3.3%	2.9%
PF09285	96%	4	8585	9.1%	4.9%
PF17136	97%	4	7896	5.4%	4.1%

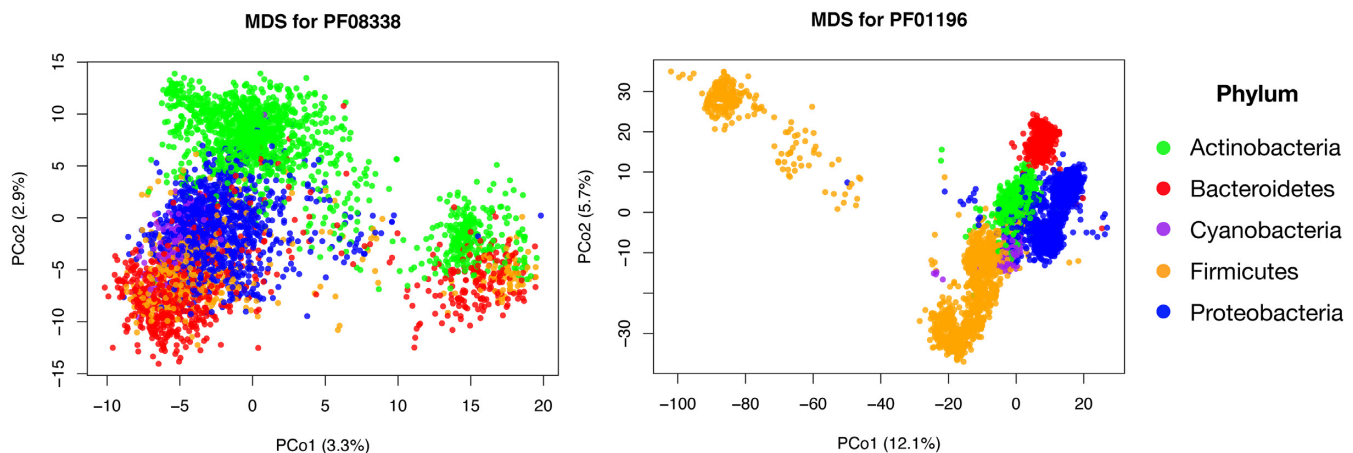


Figure 3. Phylum-based clustering for PF08338 and PF01196. For MDS analysis, only the sequences associated to the five most-abundant phyla (Proteobacteria, Firmicutes, Actinobacteria, Bacteroidetes, Cyanobacteria) have been considered.

and it had been shown to have higher precision and speed over other programs (27).

We also used the medium and large complexity mock bacterial communities from the Critical Assessment of Metagenome Interpretation (CAMI) project (45) to compare the different shotgun classification methods. CAMI is a recent community-driven initiative designed to evaluate the performance of metagenomic programs by organizing benchmarking challenges on complex and realistic microbial data. In particular, in this work we compared the taxa classification performance of Core-Kaiju, standard Kaiju and Kraken2 on the high complexity datasets of the first CAMI challenge (see <https://data.cami-challenge.org/participate>), consisting of five Illumina HiSeq microbial samples of 15 Gbp each with small insert sizes (45).

The Core-Kaiju protocol

After defining the core PFAMs, we created two protein databases for Kaiju: the first database only contains the protein sequences from the core families, whereas the second database is the standard Kaiju database based on the bacterial subset of the NCBI NR database. The protocol then follows these steps:

1. Classify the reads with Kaiju using the database with the core protein domains.
2. Classify the reads with Kaiju using the NR database to get the preliminary relative abundances for each genus.
3. Discard from the list of genera detected in (2) those having absolute abundance of less than or equal to 20 reads in the list obtained in point (1). This threshold represents our confidence level on the sequencing pipeline (see below).
4. Re-normalize the abundances of the genera obtained in point (3).

RESULTS

Comparison between methods, small mock community dataset

We evaluated the performance of both shotgun and 16S pipelines for the taxonomic classification of the three mock communities. In the top panels of Figure 4, we show the true relative genus abundance composition of the three small mock communities versus the ones predicted via the different tested taxonomic pipelines.

We then applied the Core-Kaiju pipeline to detect the biodiversity composition of the same three mock communities.

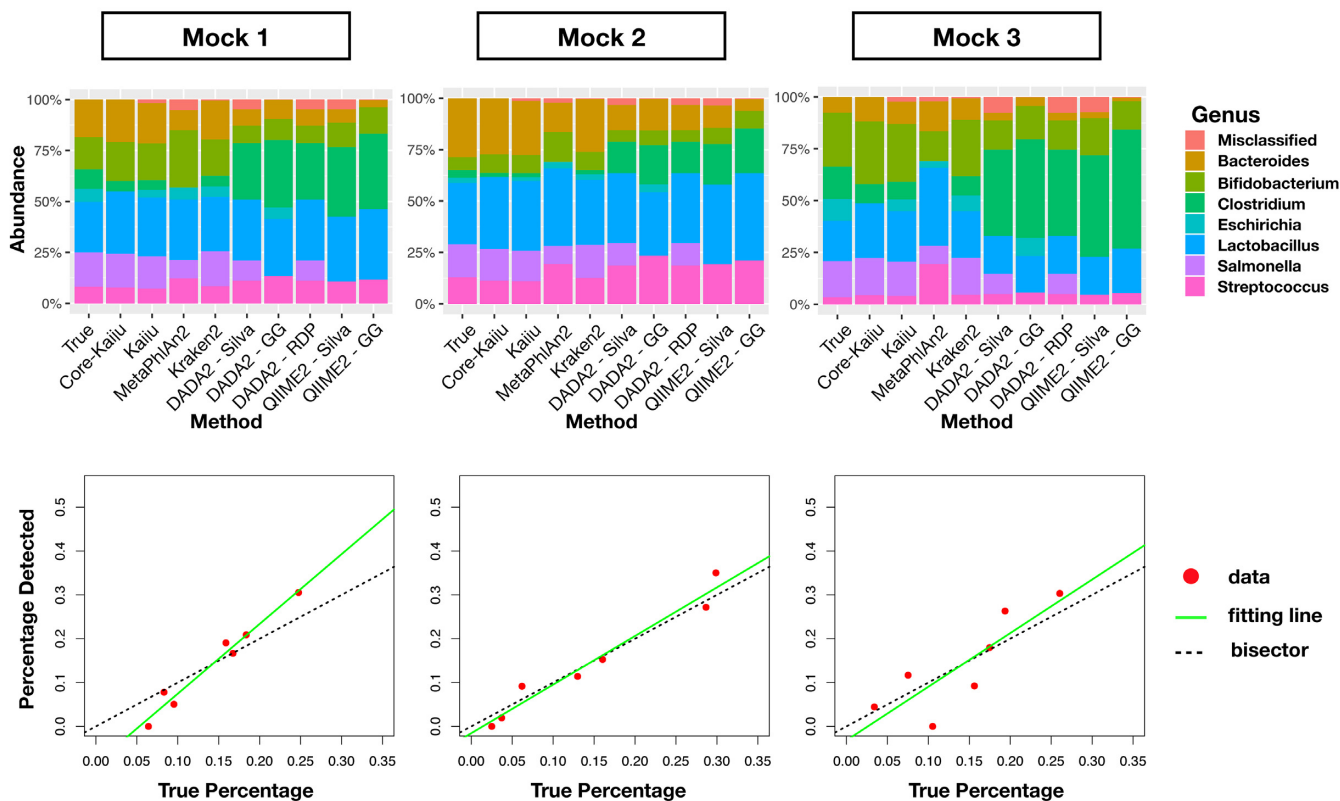


Figure 4. Comparison between theoretical and predicted relative abundances in small mock communities. Top panels: predicted relative abundance composition of the three small mock communities via different taxonomic classification methods. Bottom panels: red points represent data of relative abundance predicted for the genus level by Core-Kaiju on the three mock communities versus the true ones, known *a priori*. The green line is the linear fit performed on obtained points which, in the best scenario, should coincide with the quadrant bisector (dotted black line). In all three cases the predicted community composition was satisfactorily captured by our method, with an R -squared value of 0.97, 0.96 and 0.71, respectively.

In Figure 4, bottom panels, we plot the linear fit performed on predicted relative abundances via Core-Kaiju versus theoretical ones, known *a priori*. As we can see, in all three cases the predicted community composition was satisfactorily captured by our method, with an R^2 value higher than 0.7.

Our goal was to quantitatively compare the performance of different methods in terms of both biodiversity and relative abundances. As for the first, we chose to measure it via the F_1 score applied at the genera level. More precisely, we define the *recall* of a given taxonomic classification method as the number of truly-positive detected genera (present in a community and thus correctly detected by the method), T_p , over the sum between T_p and F_n , the number of false-negative genera (present in a community, but missed to be classified). In contrast, we define the *precision* to be the ratio between T_p and the sum of T_p and F_p , the number of false-positive genera (not present in a community and thus incorrectly detected as present). Finally, the F_1 biodiversity score is twice the ratio between the product of recall and precision and their sum, i.e. $F_1 = 2T_p / (2T_p + F_n + F_p)$. F_1 score values obtained via the different methods for the three analysed mock communities are presented in Table 3. While F_1 describes the overall accuracy in detecting the correct number of genera in the sample, R^2 gives

the correlation between the taxa abundance measured by the pipeline and the real composition of the microbial sample. Finally, we also indicated the number of genera each method predicts, \hat{G} .

Table 3 summarizes the results of the analysis, together with the R -squared values, R^2 , obtained for the linear fit performed between true and predicted relative abundances. As we can see, both Core-Kaiju and MetaPhlAn2 gave a good estimate of the number of genera in the communities (which is equal to seven), whereas all 16S methods slightly overestimated it. Moreover, both standard Kaiju and Kraken2 predicted a number of genera much higher than the true one. Finally, fit with standard Kaiju and Core-Kaiju of the predicted abundances displayed a higher determination coefficient with respect to all other pipelines, with the exception of Kraken2, which gave comparable values. However, if we focus on the F_1 score, we can notice that Core-Kaiju outperformed all the other methods in terms of precision and recall. In particular, since the pipeline led to zero false-positive and only one false negative genus (*E. coli* in all three communities), the resulting precision and recall were 1 and 0.86 for all the sampled mocks. With Core-Kaiju, we were therefore able to produce a reliable estimate of both the number of genera within the communities and their relative abundances.

Table 3. F_1 score, R -squared values and number of predicted genera. For all three analysed mock communities, we inserted the F_1 score (twice the ratio between the product of recall and precision and their sum), the R^2 value of the linear fit performed between estimated and true abundances together with the number of predicted genera, \hat{G} , with various taxonomic methods. The true number of genera is $G = 7$ for each community

		Mock 1 ($G = 7$)			Mock 2 ($G = 7$)			Mock 3 ($G = 7$)		
		F_1 score	R^2	\hat{G}	F_1 score	R^2	\hat{G}	F_1 score	R^2	\hat{G}
Shotgun	Core-Kaiju	0.92	0.97	6	0.92	0.96	6	0.92	0.71	6
	Standard Kaiju	0.02	0.97	674	0.03	0.98	501	0.02	0.94	738
	MetaPhlan2	0.86	0.46	7	0.86	0.60	7	0.86	0.08	7
	Kraken2	0.04	0.98	333	0.05	0.99	266	0.04	0.96	378
16S	DADA2 + SILVA	0.48	0.59	18	0.41	0.73	22	0.6	0.41	13
	DADA2 + GG	0.5	0.45	17	0.43	0.60	21	0.63	0.35	12
	DADA2 + RDP	0.48	0.59	18	0.4	0.73	23	0.6	0.41	13
	QIIME2 + SILVA	0.21	0.50	41	0.21	0.59	41	0.21	0.43	41
	QIIME2 + GG	0.26	0.46	32	0.26	0.50	32	0.25	0.36	33

Relative abundance vs absolute abundance thresholds

As stated in the introduction and observed above, metagenomic classification methods, such as Kaiju, often give a high number of false-positive predictions. In principle, one could set an arbitrary threshold on the detected relative abundances, for example 0.1% or 1%, to filter out low-abundance taxa that are likely false-positives. However, different choices of the threshold typically lead to very different results. The top panels of Figure 5 shows the empirical taxa abundance distribution of the 674 genera detected by Kaiju in the first small mock community. Such biodiversity number would decrease to 34, 9 or 7 if one considers only genera accounting for more than 0.01%, 0.1% and 1% of the total number of sample reads, respectively. Moreover, looking at the empirical pattern, one can notice the main gap between genera covering a fraction of $<5 \times 10^{-3}$ with respect to the total number of reads (black points) and those covering a fraction higher than 2×10^{-2} (green points), which corresponds to the genera actually present in the artificial community. One could therefore hope that, whenever such a gap is detected in the taxa abundance distribution, this corresponds to the one between false-positive and truly present taxa. However, as will be clear in the following section, this is not the case and it is not possible to set a relative threshold for the shotgun methods that works for all the mock communities.

Application to CAMI challenge dataset

We tested and compared standard Kaiju, Kraken 2 and Core-Kaiju also on medium and high complexity mock bacterial communities obtained from the 1st CAMI challenge (45), in terms of biodiversity (recall, precision, F_1 score, \hat{G}) and abundance composition (linear fit R -squared). In Table 4, we show the results for samples 1 and 5 of the high-complexity dataset (see Supplementary Table S1, for the results of the other samples). As we can see, Core-Kaiju outperformed the other methods in terms of precision. Indeed, it only slightly overestimated the true number of genera of around 10 taxa in sample 1, and 20 taxa in sample 5 (see Table 4), which is two order of magnitude lower with respect to the other methods (that predicted >1600 of taxa). On the other hand, as also shown from the bottom panels of Figure 5, when using in standard Kaiju (or Kraken 2) a relative threshold of 1% so to reduce the number of false-

positive taxa, as suggested by the previous analysis on the small mock community, the number of predicted taxa is in this case ~ 30 , therefore strongly underestimating the real biodiversity of the samples.

As for the recall, the performance of Core-Kaiju (values around 77%) stands between standard Kaiju (values around 96%) and Kraken2 (values around 65%). The combination of recall and precision led to an F_1 score around 74%, much higher than the other two pipelines (13%). Finally, as shown in Figure 6, Core-Kaiju gave also a very good estimation of the microbial composition, with an R -squared for the fit between theoretical and predicted relative abundances above 0.88, value comparable to standard Kaiju and much higher than the one obtained with Kraken2 (0.45). In Supplementary Tables S1 to S4 of the Supporting Information, we present all the results for the other high-complexity samples as well as the analyses performed on the medium-complexity challenge dataset and the sensitivity of the classification on the absolute thresholds.

Application to human gut microbiome

Last, we applied Core-Kaiju to an empirical dataset. We analysed a cohort of 26 healthy human fecal samples from the study (56) (metagenomic sequencing data are publicly available at the NCBI SRA under accession number SRP057027). We applied standard Kaiju and found on average (over the 26 samples) 2108 bacterial genera. Similar overestimation of the number of taxa of Kaiju 1.0 would be obtained also with Kraken 2, highlighting the above mentioned problem of setting the correct threshold in order to have a realistic estimation of the sample biodiversity.

The right panel of Figure 7 shows the empirical taxa abundance distribution of one individual (sample ID: SRR2145359). As we can see, in this case the only apparent gap occurs between relative abundance of $<10^{-1}$ and those >0.5 , with only one genus. It therefore results quite unrealistic that all the taxa but one should be considered false-positive. The same plot shows the vertical lines corresponding to threshold on relative population of 0.01%, 0.1% and 1% above which we have 97, 32 and 10 taxa, respectively.

In contrast, with Core-Kaiju we did not need to tune a relative threshold. Instead, by removing false-positive through the (fixed) absolute abundance of 20 reads we ended up with 21 genera (orange diamonds in Figure 7), which is compatible with previous estimates. In fact, the

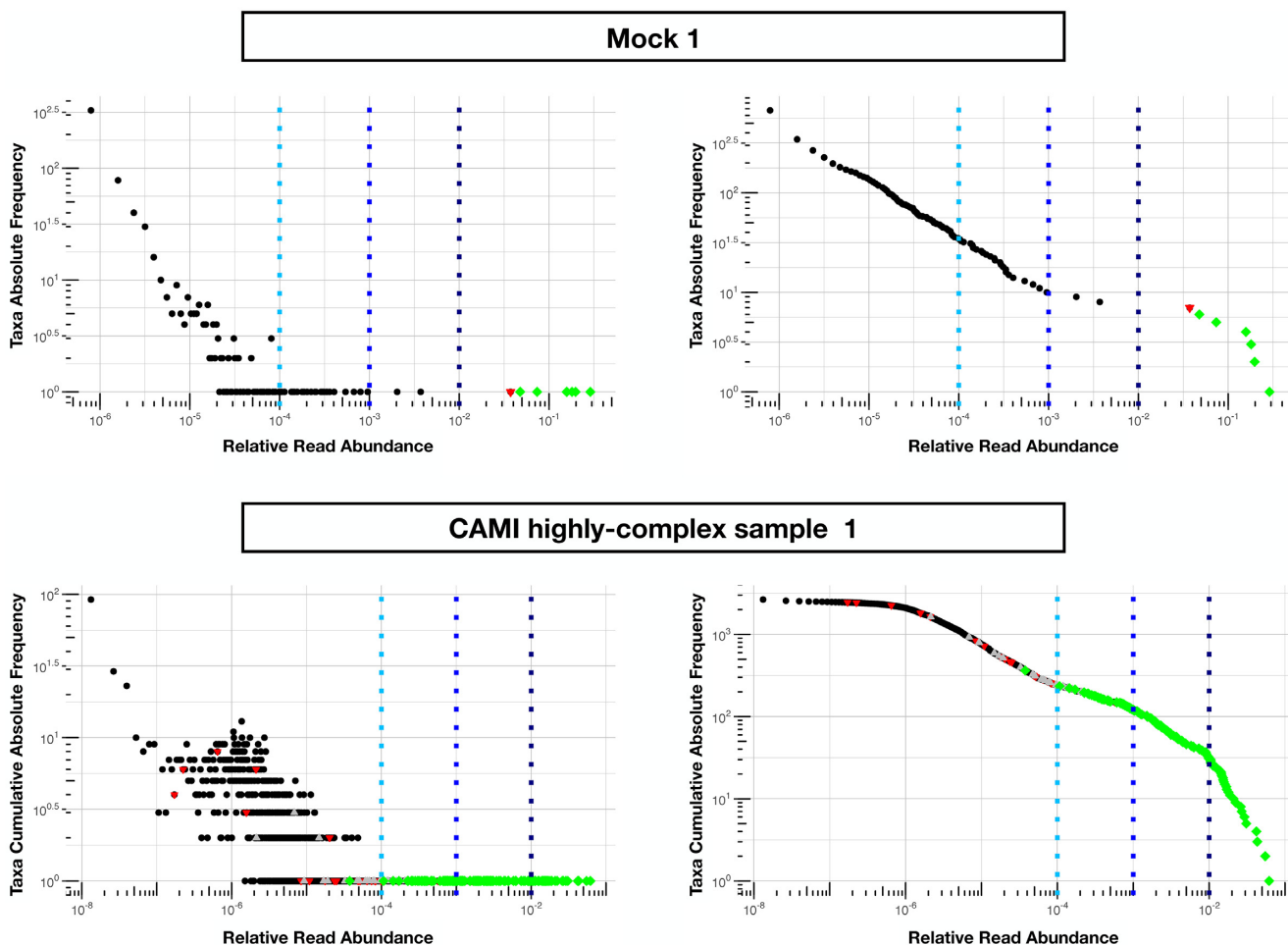


Figure 5. Relative vs absolute abundance thresholds for false-positive detection. Top panels: taxa abundance distribution plots for the first mock community (see Materials and Methods section). Green diamonds are the genera actually present in the artificial community and correctly detected by Core-Kaiju algorithm. The red triangle corresponds to the unique false-negative genus (*E.coli*) undetected with the newly proposed method. Dashed lines represent relative abundance thresholds on standard Kaiju output of 0.01%, 0.1% and 1%, respectively, which would have led to a biodiversity estimate of 34, 9 and 7 genera, respectively. Imposing an absolute abundance threshold of twenty reads on standard Kaiju output directly, would instead lead to an overestimation of 99 genera. Bottom panels: the same analyses have been performed on the CAMI high-complex sample 1. Again, green diamonds represent the 146 out of 193 genera present in the community and correctly detected by our pipeline. In this case, in addition to the remaining 47 false-negative genera (red triangles) we have also the presence of 58 false-negative genera, here represented by gray triangles. Setting a threshold on the relative abundance of reads produced by standard Kaiju gives a number of genera of 237 for the 0.01%, 120 for the 0.1% and 30 for the 1% threshold, respectively. Left and right panels represent, respectively, log-log absolute frequency and cumulative patterns of the taxa abundances in the mock communities.

available amplicon-sequencing datasets from stool samples of healthy participants of the human microbiome project (1) suggest that there are on average 25 different bacterial genera per sample (based on 174 samples with at least >5k reads per sample using 97% OTU clustering). However, in terms of taxa composition, Core-Kaiju predicted abundances are different from those obtained using 16s classification methods (1).

DISCUSSION

An important source of errors in the performance of any algorithm working on shotgun data is the high level of plasticity of bacterial genomes, due to widespread horizontal transfer (41,57-61). Indeed, most highly abundant gene families are shared and exchanged across genera, making them both a confounding factor and a computational bur-

den for algorithms attempting to extract species presence and abundance information. Thus, while having access to the sequences from the whole metagenome is very useful for functional characterization, restriction to a smaller set of families may be a very good idea when the goal is to identify the species taxa and their abundance.

To summarize, we have presented a novel method for the taxonomic classification of microbial communities which exploits the peculiar advantages of both whole-genome and 16S rRNA pipelines. Indeed, while the first approaches are recognized to better estimate the relative taxa composition of samples, the second are much more reliable in predicting the true biodiversity of a community, since the comparison between taxa-specific hyper-variable regions of bacterial 16S ribosomal gene and comprehensive reference databases allows in general to avoid the phenomenon of false-positive taxa detection. Indeed, the identification of a threshold in

Table 4. Performance comparison on CAMI high-complexity samples 1 and 5. In the first four columns, we inserted the values for the precision, the recall, the F_1 score, the R^2 value of the linear fit performed between estimated and true abundances, and the number of predicted genera \hat{G} with Core-Kaiju, standard Kaiju and Kraken2. The true number of genera is $G = 193$ for each sample. In the last column we also inserted the number of genera one would predict with standard Kaiju and Kraken2 by setting a relative threshold of 1%, i.e. by considering false-positive all those genera having a relative abundance of <0.01 in the sample. We denoted this quantity by $\hat{G}_{1\%}$.

	Sample 1 ($G = 193$)						Sample 5 ($G = 193$)					
	Precision	Recall	F_1 score	R^2	\hat{G}	$\hat{G}_{1\%}$	Precision	Recall	F_1 score	R^2	\hat{G}	$\hat{G}_{1\%}$
Core-Kaiju	0.72	0.76	0.74	0.90	204	—	0.72	0.79	0.75	0.88	213	—
standard Kaiju	0.07	0.96	0.13	0.92	2652	30	0.07	0.96	0.13	0.89	2660	26
Kraken2	0.07	0.65	0.13	0.45	1715	27	0.07	0.65	0.13	0.45	1697	26

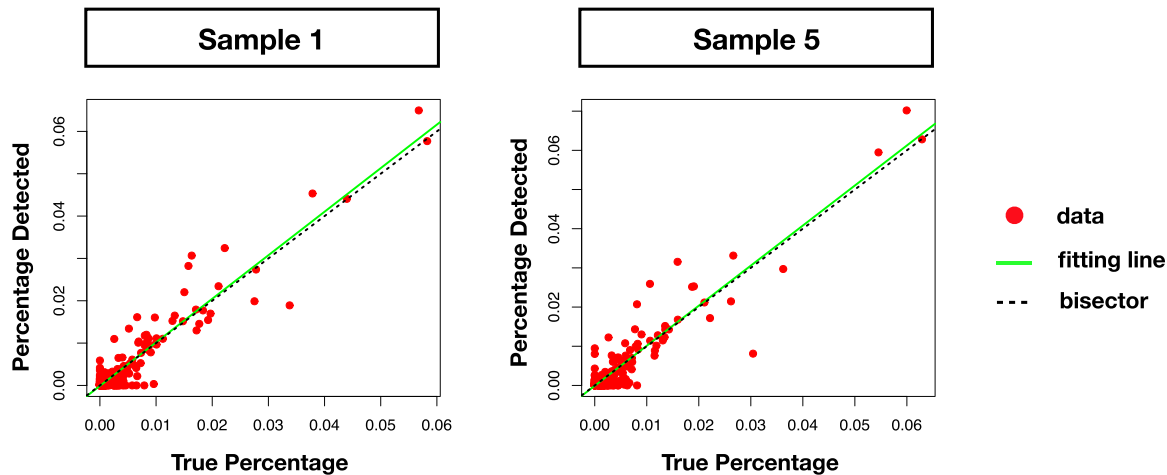


Figure 6. Linear fit between theoretical and predicted relative abundances with Core-Kaiju. Red points represent data of relative abundance predicted for the genus level by Core-Kaiju on sample 1 and 5 from the CAMI highly-complex dataset versus the ground-truth abundances, known *a priori*. The green line is the linear fit performed on such values which, in the case of perfect matching between data and Core-Kaiju output, should coincide with the quadrant bisector (dotted black line). In both cases, the predicted community composition was satisfactorily captured by our method, with a correlation with the real taxa abundances of $R^2 = 0.9$ and $R^2 = 0.88$ for sample 1 and 5, respectively.

shotgun methods to remove most of the false-positive is of course a critical problem, because in general the true taxa composition is not known, and thus setting the wrong threshold may lead to a huge over- (or under-) estimation of the sample biodiversity, as shown in this work.

Inspired by the role of 16S gene as a taxonomic fingerprint and by the knowledge that proteins are more conserved than DNA sequences, we proposed an updated version of Kaiju, an open-source program for the taxonomic classification of whole-genome high-throughput sequencing reads where sample metagenomic DNA sequences are firstly converted into amino-acid sequences and then compared to microbial protein reference databases. We identified a class of ten domains, here denoted by core PFAMs, which, analogously to 16S rRNA gene, on one hand are present in the overwhelming majority of proteomes, therefore covering the whole domain of known bacteria, and which on the other hand occur just few times in each of them, thus allowing for the creation of a novel reference database where a fast research can be performed between sample reads and PFAMs amino-acid sequences. Tested against mock microbial communities, of different level of complexity, generated in other studies (27,45) and available online, the proposed updated version of Kaiju, Core-Kaiju, outperformed popular 16S rRNA and shotgun methods for

taxonomic classification in the estimation of both the total biodiversity and taxa relative abundance distribution. In fact, by fixing an absolute threshold with Core-Kaiju (by only considering abundances greater to twenty reads), we are able to correctly classify the biodiversity in all samples of different size and complexity, while keeping a very good performance in the prediction of taxa abundances.

We highlight that other microbiome sequencing approaches exist beyond metagenomics or 16S amplicons on a MiSeq (integrated instrument performing clonal amplification and sequencing), for example PacBio long-read sequencing (62). Earl and collaborators (63) used a CAMI dataset to test the accuracy of this method and it is therefore possible to indirectly compare Core-Kaiju with PacBio through their results. Also in this case we found that our method gives a slightly higher R^2 score for the genera abundances composition, confirming the competitiveness of Core-Kaiju even with the long-read sequencing technologies, such as PacBio. However, a deeper comparison with these methods goes beyond the scope this work.

Our promising results pave the way for the application of the newly proposed pipeline in the field of microbiota-host interactions, a rich and open research field which has recently attracted the attention of the scientific world due to the hypothesized connection between human microbiome

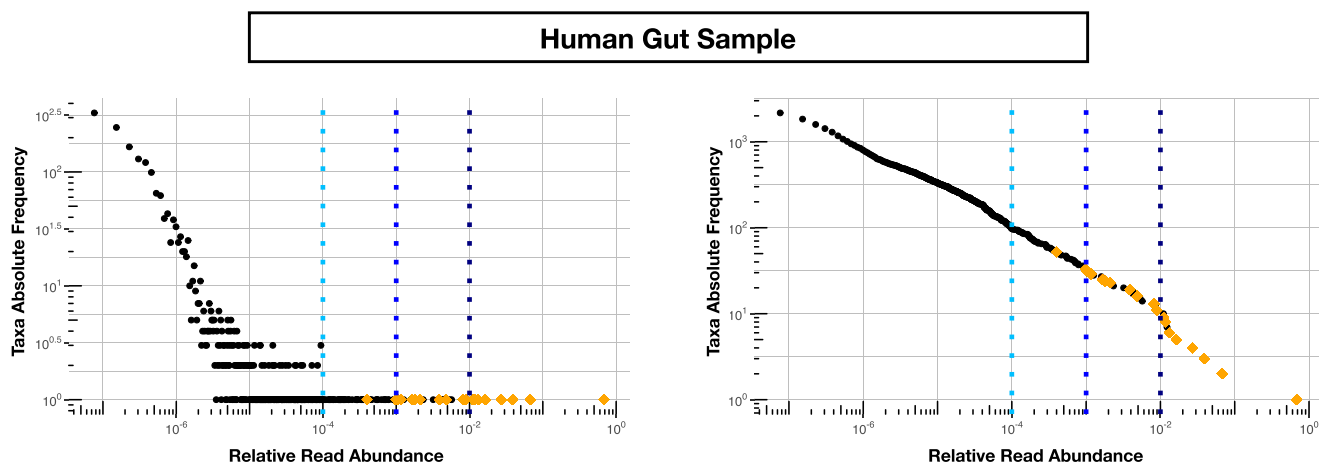


Figure 7. Relative vs absolute abundance thresholds in the human gut sample. Taxa abundance distribution plots for a human gut sample of a healthy individual, where standard Kaiju detects (without any threshold) 2165 genera. In this case the number (and label) of the actual present genera is unknown. Nevertheless estimates from a reference cohort of stool microbiomes (see <https://www.hmpdacc.org/hmp/HM16STR/>) from 174 healthy HMP participants (16S V3–V5 region, >5k reads per sample, 97% OTU clustering), report an average number of genera per sample of 25 (max = 46, min = 9) (1). Setting a threshold on the relative abundance of reads produced by standard Kaiju gives a number of genera of 97 for the 0.01%, 32 for the 0.1% and 10 for the 1% threshold, respectively. In contrast, considering false-positive all genera with less or equal to twenty reads in standard Kaiju output, we end up with 625 genera. Orange diamonds in plot correspond to the 21 genera detected with Core-Kaiju, a number compatible with the reported estimates. Left and right panels represent log–log absolute frequency and cumulative patterns, respectively.

and healthy/disease (64,65). Having a trustable tool for the detection of microbial biodiversity, as measured by the number of genera and their abundances, could have a fundamental impact in our knowledge of human microbial communities and could therefore lay the foundations for the identification of the main ecological properties modulating the healthy or ill status of an individual, which, in turn, could be of great help in preventing and treating diseases on the basis of the observed patterns.

DATA AVAILABILITY

All data and codes used for this study are available online or upon request to the authors. Raw data for the three in-silico mock communities (27) are publicly available at the Sequence Read Archive (SRA) portal of NCBI under accession number SRP059928. Metagenomic sequencing data of the healthy human fecal samples from the study (56) are publicly available at the NCBI SRA under accession number SRP057027. CAMI medium and high complexity datasets are available at <https://data.cami-challenge.org/participate> under request. All codes and scripts are freely available at <https://github.com/liphlab/Kaiju-core>

SUPPLEMENTARY DATA

[Supplementary Data](#) are available at NAR Online.

ACKNOWLEDGEMENTS

M.C.L., S.S. and A.K. gratefully acknowledge Cariparo foundation Visiting Program 2018 and the support of NVIDIA Corporation with the donation of the Titan Xp GPU also used in part for this research.

FUNDING

STARS GRANT UNIPD 2018 BioReACT (to S.S.). Funding for open access charge: UNIPD DFA DOR Funding. *Conflict of interest statement.* None declared.

REFERENCES

- Turnbaugh,P.J., Ley,R.E., Hamady,M., Fraser-Liggett,C.M., Knight,R. and Gordon,J.I. (2007) The human microbiome project. *Nature*, **449**, 804–810.
- Gevers,D., Knight,R., Petrosino,J.F., Huang,K., McGuire,A.L., Birren,B.W., Nelson,K.E., White,O., Methé,B.A., Huttenhower,C. *et al.* (2012) The Human Microbiome Project: a community resource for the healthy human microbiome. *PLoS Biol.*, **10**, e1001377.
- Thompson,L.R., Sanders,J.G., McDonald,D., Amir,A., Ladau,J., Locey,K.J., Prill,R.J., Tripathi,A., Gibbons,S.M., Ackermann,G. *et al.* (2017) A communal catalogue reveals Earth’s multiscale microbial diversity. *Nature*, **551**, 457–463.
- Bork,P., Bowler,C., De Vargas,C., Gorsky,G., Karsenti,E. and Wincker,P. (2015) Tara Oceans studies plankton at planetary scale. *Science*, **348**, 873.
- Alberti,A., Poulain,J., Engelen,S., Labadie,K., Romac,S., Ferrera,I., Albini,G., Aury,J.-M., Belser,C., Bertrand,A. *et al.* (2017) Viral to metazoan marine plankton nucleotide sequences from the Tara Oceans expedition. *Sci. Data*, **4**, 170093.
- Goldford,J.E., Lu,N., Bajić,D., Estrela,S., Tikhonov,M., Sanchez-Gorostiaga,A., Segrè,D., Mehta,P. and Sanchez,A. (2018) Emergent simplicity in microbial community assembly. *Science*, **361**, 469–474.
- Costello,E.K., Stagaman,K., Dethlefsen,L., Bohannan,B.J. and Relman,D.A. (2012) The application of ecological theory toward an understanding of the human microbiome. *Science*, **336**, 1255–1262.
- Bashan,A., Gibson,T.E., Friedman,J., Carey,V.J., Weiss,S.T., Hohmann,E.L. and Liu,Y. (2016) Universality of human microbial dynamics. *Nature*, **534**, 259–262.
- Gilbert,J.A. and Lynch,S.V. (2019) Community ecology as a framework for human microbiome research. *Nat. Med.*, **25**, 884–889.
- The Integrative HMP (iHMP) Research Network, Consortium. (2019) The Integrative Human Microbiome Project. *Nature*, **569**, 641–648.
- Lynch,S.V. and Pedersen,O. (2016) The human intestinal microbiome in health and disease. *N. Engl. J. Med.*, **375**, 2369–2379.

12. Wang, Y. and Kasper, L.H. (2014) The role of microbiome in central nervous system disorders. *Brain Behav. Immun.*, **38**, 1–12.
13. Huttenhower, C., Gevers, D., Knight, R., Abubucker, S., Badger, J.H., Chinwalla, A.T., Creasy, H.H., Earl, A.M., FitzGerald, M.G., Fulton, R.S. *et al.* (2012) Structure, function and diversity of the healthy human microbiome. *Nature*, **486**, 207–214.
14. Venter, J.C., Adams, M.D., Sutton, G.G., Kerlavage, A.R., Smith, H.O. and Hunkapiller, M. (1998) Shotgun sequencing of the human genome. *Science*, **280**, 1540–1542.
15. Hamady, M. and Knight, R. (2009) Microbial community profiling for human microbiome projects: tools, techniques, and challenges. *Genome Res.*, **19**, 1141–1152.
16. Segata, N., Börnigen, D., Morgan, X.C. and Huttenhower, C. (2013) PhyloPhlAn is a new method for improved phylogenetic and taxonomic placement of microbes. *Nat. Commun.*, **4**, 2304.
17. Tessler, M., Neumann, J.S., Afshinnekoo, E., Pineda, M., Hersch, R., Velho, L.F.M., Segovia, B.T., Lansac-Toha, F.A., Lemke, M., DeSalle, R. *et al.* (2017) Large-scale differences in microbial biodiversity discovery between 16S amplicon and shotgun sequencing. *Sci. Rep.*, **7**, 6589.
18. Langille, M.G., Zaneveld, J., Caporaso, J.G., McDonald, D., Knights, D., Reyes, J.A., Clemente, J.C., Burkepile, D.E., Vega Thurberg, R.L., Knight, R. *et al.* (2013) Predictive functional profiling of microbial communities using 16S rRNA marker gene sequences. *Nat. Biotechnol.*, **31**, 814–821.
19. Klindworth, A., Pruesse, E., Schweer, T., Peplies, J., Quast, C., Horn, M. and Glöckner, F.O. (2013) Evaluation of general 16S ribosomal RNA gene PCR primers for classical and next-generation sequencing-based diversity studies. *Nucleic Acids Res.*, **41**, e1.
20. Kembel, S.W., Wu, M., Eisen, J.A. and Green, J.L. (2012) Incorporating 16S gene copy number information improves estimates of microbial diversity and abundance. *PLoS Comput. Biol.*, **8**, e1002743.
21. Vandeputte, D., Kathagen, G., D’hoel, K., Vieira-Silva, S., Valles-Colomer, M., Sabino, J., Wang, J., Tito, R.Y., De Commer, L., Darzi, Y. *et al.* (2017) Quantitative microbiome profiling links gut community variation to microbial load. *Nature*, **551**, 507–511.
22. Angly, F.E., Dennis, P.G., Skarshewski, A., Vanwongterghem, I., Hugenholtz, P. and Tyson, G.W. (2014) CopyRighter: a rapid tool for improving the accuracy of microbial community profiles through lineage-specific gene copy number correction. *Microbiome*, **2**, 11.
23. Hugenholtz, P. and Tyson, G.W. (2008) Microbiology: metagenomics. *Nature*, **455**, 481–483.
24. Peabody, M.A., Van Rossum, T., Lo, R. and Brinkman, F.S. (2015) Evaluation of shotgun metagenomics sequence classification methods using in silico and in vitro simulated communities. *BMC Bioinform.*, **16**, 362.
25. Wood, D.E. and Salzberg, S.L. (2014) Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome Biol.*, **15**, R46.
26. Wood, D.E., Lu, J. and Langmead, B. (2019) Improved metagenomic analysis with Kraken 2. *Genome Biol.*, **20**, 257.
27. Jovel, J., Patterson, J., Wang, W., Hotte, N., O’Keefe, S., Mitchel, T., Perry, T., Kao, D., Mason, A.L., Madsen, K.L. *et al.* (2016) Characterization of the gut microbiome using 16S or shotgun metagenomics. *Front. Microbiol.*, **7**, 459.
28. Menzel, P., Ng, K.L. and Krogh, A. (2016) Fast and sensitive taxonomic classification for metagenomics with Kaiju. *Nat. Commun.*, **7**, 11257.
29. Logares, R., Sunagawa, S., Salazar, G., Cornejo-Castillo, F.M., Ferrera, I., Sarmiento, H., Hingamp, P., Ogata, H., de Vargas, C., Lima-Mendez, G. *et al.* (2014) Metagenomic 16S rDNA Illumina tags are a powerful alternative to amplicon sequencing to explore diversity and structure of microbial communities. *Environ. Microbiol.*, **16**, 2659–2671.
30. Chan, C.S., Chan, K.G., Tay, Y.L., Chua, Y.H. and Goh, K.M. (2015) Diversity of thermophiles in a Malaysian hot spring determined using 16S rRNA and shotgun metagenome sequencing. *Front. Microbiol.*, **6**, 177.
31. Poretzky, R., Rodriguez-R, L.M., Luo, C., Tsementzi, D. and Konstantinidis, K.T. (2014) Strengths and limitations of 16S rRNA gene amplicon sequencing in revealing temporal microbial community dynamics. *PLoS ONE*, **9**, e93827.
32. Quast, C., Pruesse, E., Yilmaz, P., Gerken, J., Schweer, T., Yarza, P., Peplies, J. and Glöckner, F.O. (2013) The SILVA ribosomal RNA gene database project: improved data processing and web-based tools. *Nucleic Acids Res.*, **41**, D590–D596.
33. Yilmaz, P., Parfrey, L.W., Yarza, P., Gerken, J., Pruesse, E., Quast, C., Schweer, T., Peplies, J., Ludwig, W. and Glöckner, F.O. (2013) The SILVA and ‘all-species living tree project (LTP)’ taxonomic frameworks. *Nucleic Acids Res.*, **42**, D643–D648.
34. DeSantis, T.Z., Hugenholtz, P., Larsen, N., Rojas, M., Brodie, E.L., Keller, K., Huber, T., Dalevi, D., Hu, P. and Andersen, G.L. (2006) Greengenes, a chimera-checked 16S rRNA gene database and workbench compatible with ARB. *Appl. Environ. Microbiol.*, **72**, 5069–5072.
35. McDonald, D., Price, M.N., Goodrich, J., Nawrocki, E.P., DeSantis, T.Z., Probst, A., Andersen, G.L., Knight, R. and Hugenholtz, P. (2012) An improved Greengenes taxonomy with explicit ranks for ecological and evolutionary analyses of bacteria and archaea. *ISME J.*, **6**, 610.
36. Cole, J.R., Wang, Q., Fish, J.A., Chai, B., McGarrell, D.M., Sun, Y., Brown, C.T., Porras-Alfaro, A., Kuske, C.R. and Tiedje, J.M. (2013) Ribosomal Database Project: data and tools for high throughput rRNA analysis. *Nucleic Acids Res.*, **42**, D633–D642.
37. Truong, D.T., Franzosa, E.A., Tickle, T.L., Scholz, M., Weingart, G., Pasolli, E., Tett, A., Huttenhower, C. and Segata, N. (2015) MetaPhlAn2 for enhanced metagenomic taxonomic profiling. *Nat. Methods*, **12**, 902–903.
38. Callahan, B.J., McMurdie, P.J., Rosen, M.J., Han, A.W., Johnson, A.J.A. and Holmes, S.P. (2016) DADA2: high-resolution sample inference from Illumina amplicon data. *Nat. Methods*, **13**, 581.
39. Bolyen, E., Rideout, J.R., Dillon, M.R., Bokulich, N.A., Abnet, C.C., Al-Ghalith, G.A., Alexander, H., Alm, E.J., Arumugam, M., Asnicar, F. *et al.* (2019) Reproducible, interactive, scalable and extensible microbiome data science using QIIME 2. *Nat. Biotechnol.*, **37**, 852–857.
40. Grilli, J., Bassetti, B., Maslov, S. and Cosentino Lagomarsino, M. (2012) Joint scaling laws in functional and evolutionary categories in prokaryotic genomes. *Nucleic Acids Res.*, **40**, 530–540.
41. Grilli, J., Romano, M., Bassetti, F. and Cosentino Lagomarsino, M. (2014) Cross-species gene-family fluctuations reveal the dynamics of horizontal transfers. *Nucleic Acids Res.*, **42**, 6850–6860.
42. De Lazzari, E., Grilli, J., Maslov, S. and Cosentino Lagomarsino, M. (2017) Family-specific scaling laws in bacterial genomes. *Nucleic Acids Res.*, **45**, 7615–7622.
43. Mazzolini, A., Gherardi, M., Caselle, M., Cosentino Lagomarsino, M. and Osella, M. (2018) Statistics of shared components in complex component systems. *Phys. Rev. X*, **8**, 021023.
44. El-Gebali, S., Mistry, J., Bateman, A., Eddy, S.R., Luciani, A., Potter, S.C., Qureshi, M., Richardson, L.J., Salazar, G.A., Smart, A. *et al.* (2019) The Pfam protein families database in 2019. *Nucleic Acids Res.*, **47**, D427–D432.
45. Sczyrba, A., Hofmann, P., Belmann, P., Koslicki, D., Janssen, S., Røge, J., Gregor, I., Majda, S., Fiedler, J., Dahms, E. *et al.* (2017) Critical assessment of metagenome interpretation—a benchmark of metagenomics software. *Nat. Methods*, **14**, 1063–1071.
46. Segata, N., Waldron, L., Ballarini, A., Narasimhan, V., Jousson, O. and Huttenhower, C. (2019) Metagenomic microbial community profiling using unique clade-specific marker genes. *Nat. Methods*, **9**, 811.
47. Amir, A., McDonald, D., Navas-Molina, J.A., Kopylova, E., Morton, J.T., Xu, Z.Z., Kightley, E.P., Thompson, L.R., Hyde, E.R., Gonzalez, A. *et al.* (2017) Deblur rapidly resolves single-nucleotide community sequence patterns. *mSystems*, **2**, e00191-16.
48. Mitra, S., Förster-Fromme, K., Damms-Machado, A., Scheurenbrand, T., Biskup, S., Huson, D.H. and Bischoff, S.C. (2013) Analysis of the intestinal microbiota using SOLiD 16S rRNA gene sequencing and SOLiD shotgun sequencing. *BMC Genom.*, **14**, S16.
49. Lapierre, P. and Gogarten, J.P. (2009) Estimating the size of the bacterial pan-genome. *TIG*, **25**, 107–110.
50. Mazzolini, A., Grilli, J., De Lazzari, E., Osella, M., Cosentino Lagomarsino, M. and Gherardi, M. (2018) Zipf and Heaps laws from dependency structures in component systems. *Phys. Rev. E*, **98**, 012315.
51. Pang, T.Y. and Maslov, S. (2013) Universal distribution of component frequencies in biological and technological systems. *Proc. Natl. Acad. Sci. U.S.A.*, **110**, 6235–6239.

52. Haegeman, B. and Weitz, J.S. (2012) A neutral theory of genome evolution and the frequency distribution of genes. *BMC Genom.*, **13**, 196.
53. Lobkovsky, A.E., Wolf, Y.I. and Koonin, E.V. (2013) Gene frequency distributions reject a neutral model of genome evolution. *Genome Biol. Evol.*, **5**, 233–242.
54. The UniProt Consortium. (2014) UniProt: a hub for protein information. *Nucleic Acids Res.*, **43**, D204–D212.
55. 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.* (2009) Introducing mothur: open-source, platform-independent, community-supported software for describing and comparing microbial communities. *Appl. Environ. Microbiol.*, **75**, 7537–7541.
56. Lewis, J.D., Chen, E.Z., Baldassano, R.N., Otley, A.R., Griffiths, A.M., Lee, D., Bittinger, K., Bailey, A., Friedman, E.S., Hoffmann, C. *et al.* (2015) Inflammation, antibiotics, and diet as environmental stressors of the gut microbiome in pediatric Crohn's disease. *Cell Host Microbe*, **18**, 489–500.
57. Koonin, E.V., Wolf, Y.I. and Puigbo, P. (2009) The phylogenetic forest and the quest for the elusive tree of life. *Cold Spring Harb. Symp. Quant. Biol.*, **74**, 205–213.
58. Puigbo, P., Wolf, Y.I. and Koonin, E.V. (2009) Search for a 'Tree of Life' in the thicket of the phylogenetic forest. *J. Biol.*, **8**, 59.
59. Puigbo, P., Wolf, Y.I. and Koonin, E.V. (2010) The tree and net components of prokaryote evolution. *Genome Biol. Evol.*, **2**, 745–756.
60. Puigbo, P., Wolf, Y.I. and Koonin, E.V. (2019) Genome-wide comparative analysis of phylogenetic trees: the prokaryotic forest of life. *Methods Mol. Biol.*, **1910**, 241–269.
61. Kislyuk, A.O., Haegeman, B., Bergman, N.H. and Weitz, J.S. (2011) Genomic fluidity: an integrative view of gene diversity within microbial populations. *BMC Genom.*, **12**, 32.
62. Rhoads, A. and Au, K.F. (2015) PacBio sequencing and its applications. *Genomic Proteomics Bioinformatics*, **13**, 278–289.
63. Earl, J.P., Adappa, N.D., Krol, J., Bhat, A.S., Balashov, S., Ehrlich, R.L., Palmer, J.N., Workman, A.D., Blasetti, M., Sen, B. *et al.* (2018) Species-level bacterial community profiling of the healthy sinonasal microbiome using Pacific Biosciences sequencing of full-length 16S rRNA genes. *Microbiome*, **6**, 190.
64. Shreiner, A.B., Kao, J.Y. and Young, V.B. (2015) The gut microbiome in health and in disease. *Curr. Opin. Gastroen.*, **31**, 69.
65. Foster, K.R., Schluter, J., Coyte, K.Z. and Rakoff-Nahoum, S. (2017) The evolution of the host microbiome as an ecosystem on a leash. *Nature*, **548**, 43–51.