

PhyloSift: phylogenetic analysis of genomes and metagenomes

Like all organisms on the planet, environmental microbes are subject to the forces of molecular evolution. Metagenomic sequencing provides a means to access the DNA sequence of uncultured microbes. By combining DNA sequencing of microbial communities with evolutionary modeling and phylogenetic analysis we might obtain new insights into microbiology and also provide a basis for practical tools such as forensic pathogen detection.

In this work we present an approach to leverage phylogenetic analysis of metagenomic sequence data to conduct several types of analysis. First, we present a method to conduct phylogeny-driven Bayesian hypothesis tests for the presence of an organism in a sample. Second, we present a means to compare community structure across a collection of many samples and develop direct associations between the abundance of certain organisms and sample metadata. Third, we apply new tools to analyze the phylogenetic diversity of microbial communities and again demonstrate how this can be associated to sample metadata.

These analyses are implemented in an open source software pipeline called PhyloSift. As a pipeline, PhyloSift incorporates several other programs including LAST, HMMER, and pplacer to automate phylogenetic analysis of protein coding and RNA sequences in metagenomic datasets generated by modern sequencing platforms (e.g. Illumina, 454).

PhyloSift: phylogenetic analysis of genomes and metagenomes

Aaron E. Darling^{1,*}, Guillaume Jospin¹, Eric Lowe¹, Frederick A. Matsen IV⁴, Holly M. Bik¹, Jonathan A. Eisen^{2,3}

1 Genome Center, University of California, Davis, California, United States of America

2 Department of Evolution and Ecology, University of California, Davis, California, United States of America

3 Department of Medical Microbiology and Immunology, University of California, Davis, California, United States of America

4 Fred Hutchinson Cancer Research Center, Seattle, Washington, United States of America

* E-mail: aarondarling@ucdavis.edu

Introduction

Metagenomics - the sequencing of DNA isolated directly from the environment - has become a routinely used tool with wide applications [57]. Used primarily in the study of microorganisms, metagenome sequencing has now been carried out on a variety of environments where one finds microbes - from plants and animals to every kind of natural and man-made environment around the globe. Metagenomic sequencing has provided fundamental insight into the diversity of microbes and their function and roles in ecosystems. Initially, metagenomics was used largely as a way of simply obtaining some genomic information about organisms for which culturing technique was unknown[4]. However, due to decreases in the cost and the difficulty of sequencing, metagenomics has become a tool for studying any microbial community, regardless of cultivability.

One strength of metagenomic approaches arises from the ability to sample the genomes of organisms in a particular environment approximately uniformly at random. This effect is achieved with the random “shotgun” sequencing methods originally applied for *de novo* genome sequencing of individual organisms [60, 59]. From random shotgun sequence data of DNA isolated from environmental samples, one can make inferences about what organisms are present in a sample (i.e., who is there?) as well as their functional potential (i.e., what are they doing?). In addition, by comparing shotgun metagenomic data across samples one can study larger scale issues such as ecology and biogeography and also attempt to correlate particular organisms or functions with “metadata” about samples (e.g., health status, nutrient cycling rates, etc [58]). Furthermore, by sampling a community directly one can avoid certain problems inherent in culturing such as contamination, population bottlenecks, and taxonomic bias [17]. In this sense metagenomics can be considered an extension of “culture-independent” ribosomal RNA gene surveys [23]. The great potential for novel insight into microbial communities has led researchers in fields as diverse as medicine and agriculture, law enforcement, biodefense, ecology, evolution, and industry to apply metagenomic methods.

Although great potential exists for metagenomics to yield insight into the hidden world of microbes, many challenges remain before this potential can be realized. Perhaps the biggest challenges lie in analysis of the data [12]. First, metagenomic samples reflect entire communities of organisms, unlike “traditional” genome sequencing of individual organisms or clones (i.e., from cultures of a single isolate where genetic diversity has undergone a bottleneck). The large number of microbial taxa in environmental samples can be a challenge for some types of analysis. Within-species genomic polymorphism presents an even greater challenge [29]. This challenge arises largely because shotgun metagenomic sequencing protocols destroy some of the most valuable information present in a sample: genetic linkage. Loss of linkage information occurs in two ways: during sample extraction and fragmentation of DNA for sequencing. In nearly all metagenomic sample processing methods, cells from the microbial community are lysed together to obtain a common pool of DNA. This practice causes DNA from many different cells to mix together, so that the cellular compartmentalization of individual genotypes is destroyed. Subsequently, long chromosome-scale DNA fragments are typically broken by mechanical or enzymatic means into fragments small enough for processing with current sequencing protocols. The resulting sequenced fragments are usually less than 1 Kbp in length. Although it is possible to generate data for larger fragments via cloning [4] or using Pacific Biosciences sequencing, most metagenomic data is currently being generated with short read/short insert sequencing chemistry such as that offered by Illumina. Though short read methods are quicker, easier, and lower cost per base and per read than large fragment approaches, there is a tradeoff in information quality. The shearing results in further loss of genetic linkage information, since we no longer have direct information on how short DNA fragments are arranged into chromosome-scale molecules.

The lack of linkage information limits the ability to use metagenomic data for phylogenetic and population genetic analysis, since most current methods assume complete linkage information is available. In practice, improved sample processing methods could potentially retain the genetic linkage information of a microbial community throughout the sequencing process. High throughput single-cell genomics (e.g. applied to hundreds or thousands of cells) offers an alternative to the standard metagenomics workflow

that preserves information about the compartmentalization of genetic material into cells [62, 32, 48]. However, single cell approaches are still limited in their utility by a number of technical issues including contamination, expensive and extensive equipment needs, missing data, and the creation of chimeras [6]; they will always be more limited in throughput than their standard metagenomics counterparts.

Thus the research community is left with developing and using computational methods to sift through and make sense of short read, random shotgun metagenomic data. Though there are many important steps in analyzing metagenomic data, we believe that a critical component is phylogenetic analysis of the sequences. Among the uses of phylogenetic analysis in metagenomics are: improved classification of sequences using phylogenetic methods, functional prediction for genes, alternative metrics of alpha and beta diversity, improved identification of operational taxonomic units (OTUs), and sequence binning [40, 37, 19, 27, 25, 54, 65, 9, 8, 52, 26, 55, 18].

In the present manuscript, we introduce PhyloSift, a new method for phylogenetic analysis of metagenomic samples and for comparison of community structure among multiple related samples. The new method leverages phylogenetic models of molecular evolution to provide high resolution detection of organisms in a metagenome. Our approach is based on well known statistical phylogenetic models, is amenable to Bayesian hypothesis testing, and uses name-independent and OTU-free analyses to provide higher resolution about microbial community assemblages (versus methods that rely on taxonomy or OTUs). These methods can be applied to any single phylogeny at a time, and expand on our previous experience building AMPHORA [66]. We additionally propose a set of 37 “elite” marker gene families that have largely congruent phylogenetic histories, thus improving the limit of detection for rare organisms in microbial communities. We contribute an open-source implementation of the method that has been engineered for ease-of-use on 64-bit Linux and Mac platforms. Finally, we compare the features of PhyloSift to some related methods to provide readers with insight into when use of our approach is and is not appropriate.

Previous work

Estimating community composition from amplicon data

High throughput sequencing of marker gene amplicons (homologous loci such as 16S/18S rRNA) has emerged as a powerful and straightforward means to analyze microbial community structure. In contrast to shotgun metagenomics, amplicon approaches currently make the detection of rare taxa easier and require less starting genomic material than some metagenomic approaches, although transposon-catalyzed libraries have been generated from as little as 30 pg total material [2]. By design rRNA surveys offer a “standardized” snapshot of microbial communities with reads from a single or small number of genes, considerably simplifying the tasks of alignment and analysis. Amplicon studies generally focus on characterizing and comparing microbial community structure without much analysis of functional gene repertoire.

A variety of software pipelines can be used to process and analyze rRNA amplicon data [5]. Inferring microbial assemblages typically relies on clustering of Operational Taxonomic Units (e.g. at a 97% sequence identity cutoff, using either *de novo* or reference-based clustering), where taxonomy is assigned to representative sequences using either BLAST searches or the RDP classifier (a Naive Bayesian Classifier [61]). Users can subsequently carry out a suite of downstream ecological and diversity analyses, including rarefaction (e.g. analyses for Chao1 estimation, OTU richness, or phylogenetic diversity as implemented in QIIME [11]), and Principal Component Analysis and Jackknife cluster analysis (e.g. using phylogeny-derived UniFrac distances [35]).

Amplicon approaches are now relatively cheap and easy to carry out. However some computational bottlenecks hinder fine-scale analysis of amplicon data. Analysis pipelines can not readily distinguish rare members of a microbial community from noise in data caused by sequencing errors or chimeric reads [5]. The RDP classifier [61] provides a statistical method for assessing confidence in taxonomic classifications.

Any of these methods are limited relative to phylogenetic methods, in that they can only distinguish named groups of organisms and are limited to the resolution of the taxonomy.

Community composition from metagenomes

Methods have also been developed to estimate and analyze community composition from metagenomic data sets. These methods typically focus on a small subset of widely conserved marker genes mined from metagenomic sequence reads, usually representing 1% of any given shotgun dataset. Marker genes include well-characterized protein coding genes (e.g. ribosomal proteins or elongation factor genes) or conserved noncoding regions (e.g. rRNA). A variety of computational approaches are now available to investigate the community composition of metagenome datasets, including: AMPHORA (bacterial protein markers and tree insertion via parsimony) [66] and AMPHORA2 (bacterial/archaeal protein and DNA markers and tree insertion via likelihood or parsimony) [65], MLTreeMap (reference gene families with taxonomic and functional information and tree insertion via maximum likelihood) [54], MetaPhyler (taxonomic classifiers for each of the reference marker genes published in the AMPHORA set) [33], EMIRGE (an expectation-maximization method to reconstruct rRNA genes from metagenome data and estimate taxon abundance) [41], and PhylOTU (phylogenetic methods to mine rRNA and define OTUs from metagenome data)[52].

An interesting alternative approach is employed by the software MetaPhlan [51], which instead of using universally conserved genes, employs a database of clade-specific genes to estimate abundance of known taxonomic groups. This approach may work well in environments where the genomic diversity is very well characterized.

Community composition analysis from metagenomes has some potential advantages over amplicon studies. For example, metagenome sequencing might avoid bias introduced by preferential binding of PCR primers to DNA from some organisms in amplicon studies and can also capture genomes from organisms which lack amplicon target genes, such as viruses. Whole-metagenome surveys also have the potential to provide insight into enzymatic and other functional processes in microbial communities, and so a single dataset can provide both community composition and functional information. One major limiting factor is that reference genome databases have narrow phylogenetic breadth relative to marker genes (e.g. rRNA) [63].

Taxonomic classification of metagenome sequences

Current methods for taxonomic classification of metagenomic sequences generally leverage one or two information sources: sequence composition and/or sequence identity to reference databases. Some existing composition classifiers include TACOA (supervised classification using k-nearest neighbors) [13], PhyloPythia [39] and PhyloPythiaS (multiclass support vector machine classifier using oligonucleotide frequencies) [46], NBC (Naive Bayesian Classifier) [49], and Eu-Detect (oligonucleotide binning to separate eukaryote sequences in feature vector space) [42], although this is not an exhaustive list. Related methods such as Self-Organizing Maps (e.g. eSOMS [14]) can be applied to tetranucleotide frequencies in combination with other information sources such as contig coverage/abundance information to produce visual "maps" displaying different bins, although this does not result in taxonomic assignment.

Identity-based classification methods compare metagenome sequences against reference databases to identify putative homologs. Examples of current identity-based classification tools include MEGAN (a Lowest Common Ancestor algorithm that summarizes BLAST outputs to assign taxonomy) [24], SORT-Items (reciprocal BLAST approach to detect significant orthology) [43], MTR (a variation on Lowest Common Ancestor approaches that uses multiple taxonomic ranks) [22], and ProViDE (analysis of alignment parameter thresholds, specifically customized for classifying viral sequences) [21]. Some approaches are able to combine both sequence identity and composition when classifying [9, 8]. Again, this is not an exhaustive list.

As the focus of our current work is on phylogenetic analysis rather than taxonomic classification, we do not discuss the relative merits of each approach to taxonomic classification in detail, nor do we provide benchmarks of taxonomic classification methods.

Methods

PhyloSift implements a method for analyzing microbial community structure directly from metagenome sequence data. Figure 1 gives an overview of the analysis workflow as executed when analyzing a metagenomic sample. The analysis can be decomposed into four stages: 1. searching input sequences for identity to a database of known reference gene families; 2. adding input sequences to a multiple alignment with reference genes; 3. placement of input sequences onto a phylogeny of reference genes; and 4. generation of taxonomic summaries. We now describe the details of each step along with our design decisions and rationale.

Reference gene families used by PhyloSift

The standard PhyloSift database includes a set of 37 “elite” gene families previously identified as nearly universal and present in single-copy. These 37 gene families are a subset of the 40 previously reported [64], with three families excluded because they frequently have partial length homologs in some lineages. These “elite” families represent about 1% of an average bacterial genome, as estimated from current genome databases. In other work we have demonstrated that phylogenetic trees reconstructed on individual genes in this set are generally congruent with each other [30, 48], suggesting that concatenating alignments of these families will yield a valid and more powerful estimate of their phylogenetic history. Other groups have also demonstrated that trees inferred from concatenate alignments demonstrate the least conflict with trees inferred separately from other microbial amino acid sequences [1]. During the database update process (described below), these gene families are automatically extended to include putative homologs from eukarya and some viruses with large genomes such as the Mimivirus. Most small viral genomes lack homologs of these gene families.

In addition to the elite 37 families, the PhyloSift database also includes four additional sets of gene families:

- 16S and 18S ribosomal RNA genes
- mitochondrial gene families
- Eukaryote-specific gene families
- Viral gene families

Combined, this yields a set of approximately 800 gene families in the standard PhyloSift database, most of which are viral.

Detailed PhyloSift client workflow

Sequence identity search

This first step in a PhyloSift analysis aims to identify regions of the input sequences that may be homologous to gene families in the reference database. Input sequences to this step can be of any length ranging from short 30nt next-generation sequence reads to fully assembled genomes or metagenomes. Recognized input formats include FastA and FastQ (paired, unpaired, phred33, phred64, and/or interleaved pairing), and these can optionally be supplied as bzip2 or gzip compressed data files. Sequence input can be streamed via stdin or unix named pipes. Amino acid input sequences can also be processed.

PhyloSift uses LAST [28] for sequence similarity search against the reference databases. We evaluated many possible search algorithms and implementations before finally selecting LAST. Other options we evaluated were BLAST [3] v2.2.23, BLAST+ [10] v2.2.28+, and RAPsearch2 [68] v2.04, and bowtie2 [31] v2.0.0-beta5. Given the large volume of sequence data that must be processed, a key evaluation criterion was algorithm efficiency both in CPU time and memory requirements. A second criterion is the ability to perform six-frame translated searches of DNA sequence against an amino acid database with the possibility to tolerate frame-shift errors in the sequence. Among the evaluated methods, BLAST and BLAST+ were slowest (data not shown) and frameshift detection was non-functional in the version of BLAST+ we obtained from NCBI. We excluded these from further consideration. RAPsearch2 was much more computationally efficient than either BLAST or BLAST+, but the version we obtained could not process sequences > 1kbp and did not support frameshift detection. In our testing, LAST was able to process sequence data as quickly as RAPsearch2 (e.g. orders of magnitude more quickly than BLAST) and supports both frameshift detection and input sequences of arbitrary length. LAST also supports all three of the primary search types we require: DNA vs. DNA, DNA vs. AA, and AA vs. AA. We also evaluated bowtie2, a program typically used for mapping reads to a reference genome, for the purpose of screening reads against a database of noncoding RNA sequences (currently 16S and 18S). bowtie2 does not offer translated amino-acid searches. Relative to LAST, bowtie2 is able to identify similarity to the RNA database sequences more quickly. However, even though the speedup over LAST was substantial (data not shown), the compute time saved is small relative to the total time consumed in the complete PhyloSift client workflow. Therefore we decided to use only LAST since using only a single local alignment search tool simplifies the code. One shortcoming of LAST is that current versions do not support multithreaded parallelism. PhyloSift implements optional process-level parallelism by spawning multiple LAST searches against the protein database.

One feature of reference gene family sequences being searched at this stage bears special mention. During database construction (described elsewhere) a representative subset of all available sequences are selected from each gene family to be searched in the search stage. These representatives are chosen to span the phylogenetic diversity of the gene family without including closely related sequences (see Section “PhyloSift database update workflow”). This is important because it reduces the volume of sequence to search and because part of LAST’s fast heuristic to identify candidate regions to align involves eliminating redundant and repetitive *k*-mers from the search space [28]. Thus, a database constructed with all sequences (and not just divergent representatives) could in principle reduce sensitivity in aligning reads to those database sequences.

The search stage identifies a set of candidate amino acid sequences from the input data that are similar to reference gene families. If DNA was provided as input the corresponding DNA sequences are also reported.

Alignment to reference multiple alignment

Prior to the alignment stage all input sequence regions with putative homology to reference gene families have been identified and extracted. In this stage, each candidate sequence is added to an amino acid or cDNA multiple sequence alignment of the reference gene family. If the input sequences were DNA, a codon multiple sequence alignment congruent to the amino acid alignment is also generated.

PhyloSift applies the `hmmalign` program from the HMMER 3.0 software package [15] to add the candidate sequences to reference multiple sequence alignments. During construction of the PhyloSift reference database (described in section “Custom gene families”) a profile-HMM is generated from a multiple alignment of the gene family reference sequences. When processing candidate sequences, PhyloSift then uses the profile-HMM to map the input sequence to the reference multiple alignment. Application of a profile-HMM to align highly divergent sequences suffers some documented shortcomings, in particular that alignment accuracy decreases with divergence of source sequences used to construct the profile-HMM [34]. This is one avenue for future improvement of PhyloSift and protein evolution models

in general.

Finally, PhyloSift concatenates the alignments of the 37 elite markers to a single multiple sequence alignment. When a single input sequence aligns to multiple genes, the aligned sequence becomes a single row in the concatenated alignment. All other sequences are represented in separate alignment rows.

PhyloSift treats input sequences with similarity to non-coding RNAs differently than protein genes. Sequences longer than 600nt are aligned using Infernal's `cmalign` program with the global alignment option. Short sequences are aligned with `hmmalign` to a profile-HMM of the non-coding RNA molecule. Although the profile-HMM does not capture secondary structure, the alignment computation is significantly faster with currently available versions of Infernal and HMMER. In our experience a banding threshold (a parameter that determines the size of the search space and hence amount of computational effort) of 1×10^{-20} is required to obtain accurate local alignments with Infernal for short sequences, but this requires several minutes of CPU time per aligned sequence, which is not practical when aligning millions of amplicon sequences.

Placement on a phylogenetic reference tree

At this stage, aligned input sequences are placed onto a phylogenetic tree of the reference sequences. PhyloSift employs pplacer [37] for this task. pplacer can be run in either maximum likelihood (ML, the default) or Bayesian mode. When run in ML mode, pplacer identifies and reports a set of most likely attachment points for each aligned sequence to the reference phylogeny, as well as a "likelihood weight ratio" representing the relative likelihood for the chosen attachment point over other possible attachment points.

When run in Bayesian mode, pplacer calculates the posterior probability that the query sequence diverged from particular branches of the reference tree via direct integration. In contrast to ML placement which selects a single most likely attachment point, the branch posterior probability integrates over all possible attachment points for the query sequence on the branch. The posterior probability is used when calculating Bayes factors for lineage tests, described below.

Taxonomic summary of read placements

At this final stage of analysis, PhyloSift summarizes the phylogenetic placements in a human-friendly format. For each gene family, the PhyloSift database contains a gene-tree/taxonomy reconciliation encoding a pre-computed mapping of edges in the gene family phylogeny to edges in the NCBI taxonomy. The method used to calculate these reconciliations is described in the database update workflow section, below.

Input to this stage of analysis is one or more "jplace" format [36] files containing an edge-labeled reference tree for a gene family along with a collection of one or more sequence placements onto that tree. Information about each sequence's placement consists of the log-likelihood of placement at several (usually up to 7, a configurable limit) of the highest likelihood edges on the reference tree, along with the probability mass that the sequence belongs at that position of the tree, and finally the weight of the sequence. When analyzing unassembled reads the sequence weights are typically always 1, when analyzing assembled contigs the weights may be set to a value based on estimated depth-of-coverage for that contig.

PhyloSift parses each of the jplace files and uses the gene-tree/taxonomy reconciliation to convert probability mass over read placements into a probability mass over the taxonomy, summing these masses over all reads and gene families. Any particular edge in the gene tree may be mapped to many equally optimal locations in the taxonomy. PhyloSift distributes the placed sequence's mass equally among all optimal locations.

Finally, PhyloSift reports the summarized taxonomy probability mass distribution in a variety of formats.

Visual presentation of taxonomic summary

For easy visualization and exploratory data analysis, PhyloSift produces Krona plots [45] showing taxonomic probability mass in the 37 elite gene families, and a separate Krona plot showing taxonomic probability mass distribution summed across the elite families and all other families.

Figure 5 provides an example of PhyloSift's Krona reports.

Parallelism and stream computing

PhyloSift supports streaming input of sequences, this permits analysis to proceed as sequences arrive over a network connection, for example.

Comparison among samples

One of the unique aspects of PhyloSift relative to other methods for comparative metagenomics is that the phylogenetic approach we have implemented enables direct comparison of the phylogenetic structure and relative abundance of metagenome samples without resorting to taxonomic relative abundance estimates. Perhaps the most powerful exploratory data analysis tool for comparing community structures among samples is Edge Principal Component Analysis, or edge PCA [25]. Edge PCA applies the standard dimensionality-reduction tool of PCA to a matrix where columns correspond to edges in the reference phylogeny, rows correspond to each sample, and each entry is the difference in placed sequence probability masses on either side of that edge. When applied in this manner, the eigenvalues of each eigenvector that results from PCA correspond to weights indicating how important each edge in the reference phylogeny is for explaining the variation among samples in that dimension. These eigenvectors can be naturally visualized as thickened branches along the reference phylogeny [25].

PhyloSift includes the guppy program from pplacer, which in addition to edge PCA also provides means for hierarchical clustering of multiple samples using an algorithm specialized to the case of masses on a tree, calculation of Kantorovich-Rubenstein distances among samples [19], and other tools for calculating sample summary statistics such as weighted phylogenetic diversities.

PhyloSift database update workflow

An integral component of PhyloSift is an automated means to update the gene family database with newly sequenced genomes. Genome databases continue to grow quickly, with, on average, dozens of new genome sequences becoming available every week. The quality of these genomes can be highly variable, ranging from low-quality drafts to nearly finished sequence. PhyloSift's database update mechanism incorporates some basic quality control measures.

Acquiring new genome data

The PhyloSift database update module maintains a local repository of all known and processed genomes. Upon initiating a new update, the database update module identifies any new genomes available in the NCBI finished, NCBI draft, NCBI WGS, and EBI viral, organelle, bacterial, archaeal, and eukaryal databases. Any new genomes are fetched and stored in the local repository.

Gene family search and alignment workflow on each genome

In this stage, the search and alignment stages of the previously described PhyloSift client workflow are run for each new genome. After this stage, the regions from each new genome that are highly similar to gene families in the database are identified, extracted, and aligned using the family's profile-HMM. A complete multiple alignment for each family is then created by adding the aligned regions from each

genome to a single multiple alignment file. Because each region has been aligned to the same profile-HMM (or covarion model for noncoding genes) and non-aligning sites in the query genome removed, generation of a new multiple alignment is a simple matter of concatenating the individual alignments.

PhyloSift also generates codon alignments for each protein-coding gene family at this stage by replacing amino acids with their codons and replacing each gap with a gap triplet.

We note that profile-HMMs are not recomputed during the database update, thereby avoiding problems with model drift.

The PhyloSift reference database is available independently of the software at the following location: http://edhar.genomecenter.ucdavis.edu/koadman/phylosift_markers

Phylogenetic inference and pruning

The next step of database update involves constructing a phylogenetic tree for each gene family. Currently PhyloSift employs FastTree 2.1 [47] to generate approximate maximum likelihood trees for this task. PhyloSift also infers separate trees for the codon and amino acid alignments of each gene family.

Reference databases frequently contain genomes for a multitude of closely related strains. In many gene families, the gene sequences present in genomes of closely related strains may be identical to each another. Identical gene sequences would create uncertainty in the placement of reads in a strain group. In order to reduce compute time and memory requirements, closely related sequences are pruned from the PhyloSift reference database. Pruning is done with an algorithm that maximizes phylogenetic diversity of the sequence set without including any sequence pairs separated by fewer than X amino acid (or nucleotide) substitutions per site, where X is a configurable variable with default value 0.01.

Selection of representatives for similarity search

The PhyloSift client workflow uses LAST to search for similarity between input sequences and reference sequences. During the database update the set of reference sequences is updated to include representatives of any newly sequenced genomes. As above, we select a subset of sequences that maximize phylogenetic diversity while requiring sequence pairs to be separated by at least X amino acid substitutions per site. In this case, X defaults to 0.1.

Taxonomic reconciliation

Many of the data sources for new genomes provide a taxonomic identifier for the genome that places it in the NCBI taxonomy. Throughout the database update process, the associations between taxon ID and individual sequences are maintained. The tips of reconstructed phylogenies can therefore have some or all nodes annotated with the taxon ID associated with that tip. Given this information, PhyloSift generates a mapping of edges (e.g. the edge above each node) in the gene tree phylogeny to edges in the taxonomic tree. To do so, we first compute the split (bipartition) encoding of the gene tree and the taxonomic tree. A tree's split encoding is simply the set of splits encoded by each edge in the tree, where the split for edge i is a binary vector $S_i = \{s_{i,1} \dots s_{i,n}\}$, $s_{i,j} \in \{0, 1\}$. Here n is the number of leaf nodes shared by the two trees. For convenience, we denote the split encoding for the gene tree as $S^{(G)}$ and use $S^{(T)}$ for the taxonomic tree. Then for each edge i in the gene tree, we compute its mapping M_i to taxonomic tree edges as:

$$M_i = \operatorname{argmin}_{S_j \in S^{(T)}} H(S_i^{(G)}, S_j)$$

Where $H(\cdot, \cdot)$ is defined as the Hamming distance among equal-length binary vectors. We note that there may be many possible edges in $S^{(T)}$ with equally minimal Hamming distance to an edge i in $S^{(G)}$. In this case M_i includes all of these edges, and so $M_i \subseteq S^{(T)}$ and $|M_i| \geq 1$. In the client workflow when assigning placement probability mass to names, the placement mass on edge $S_i^{(G)}$ is divided equally

among the taxonomic groups associated with M_i . Finally, we discard highly ambiguous mappings where $|M_i| > y$. Here y is an ad-hoc threshold with a default value of 30. These gene tree edges are labeled “Unclassifiable” due to their extreme topological discordance with the NCBI taxonomy.

Custom gene families

PhyloSift also supports the addition of custom gene families to its database. To add a gene family to the database, a multiple sequence alignment must be provided. Optionally, a table mapping each sequence identifier in the alignment and its NCBI taxon ID may also be provided. Given these inputs, PhyloSift will construct a phylogenetic tree, create a pruned set of representative sequences for similarity searching, construct a profile-HMM for alignment, and if taxon information was provided will also compute a reconciliation between the gene tree and taxonomy. The tree-building and reconciliation steps follow the approach outlined above in the PhyloSift database update workflow, with the exception that codon alignments are not generated. The resulting data is called a “package,” and is copied into the user’s PhyloSift database. The new package will be automatically included in any future runs of the PhyloSift client workflow.

Results

Bayesian hypothesis testing for the presence of phylogenetic lineages

For various applications (e.g. microbial forensics) a practitioner might want to test for the presence of a particular lineage of interest in a metagenomic sample. Phylogenetic analysis of metagenomic reads has the potential to offer resolution beyond what would be available from taxonomic methods for metagenomics. Whereas taxonomic methods can provide resolution at specific levels in the taxonomic hierarchy, such as species, genus, etc., phylogenetic methods might be able to distinguish different subtypes of named species or novel lineages at higher taxonomic levels. Phylogenetic methods are limited only by the resolution of the reference genome phylogeny and not by the resolution of manually curated taxonomies. Phylogenetic inference has the further advantage that it is based on a statistical model of sequence change where the marginal likelihood of the data given the model $P(D|M)$ is well defined, making it possible to conduct model-based hypothesis tests using phylogenies. Taxonomic analysis methods for metagenomics are frequently based on machine learning classification methods which do not always lend themselves to such hypothesis testing.

PhyloSift provides a means to conduct Bayesian hypothesis testing for the presence of one or more query sequences belonging to organisms that have diverged along specific branches of the reference phylogeny. In order to describe the Bayesian hypothesis test we introduce the following notation: assume we are given a reference phylogenetic tree T consisting of $n > 1$ branches $\{t_1 \dots t_n\}$. Further assume we are given a collection S of sequences $s_1 \dots s_m$ which are homologous to and aligned to the sequences at the leaf nodes of the reference phylogeny. We denote the marginal likelihood that a particular sequence s_j diverged along branch t_i of the reference phylogeny as $P(s_j|t_i)$. Calculation of this marginal likelihood is implemented in the pplacer software and described elsewhere [37].

The null hypothesis we wish to test is that there are no sequences diverging from a set of one or more lineages of interest $T_x \subseteq T$. We can express the marginal likelihood of the null hypothesis M_0 as:

$$P(D|M_0) = \prod_{s_j \in S} \left[1 - \sum_{t_i \in T_x} P(s_j|t_i) \right] \quad (1)$$

which can be interpreted as the product over all sequences of the probability that the sequence does not derive from a lineage of interest in T_x . The marginal likelihood of the alternative hypothesis, e.g. that

one or more reads derive from a lineage in T_x , can simply be expressed as:

$$P(D|M_1) = 1 - P(D|M_0) \quad (2)$$

Using these marginal likelihoods we can construct a Bayes factor:

$$K = \frac{P(D|M_0)}{P(D|M_1)} \quad (3)$$

The Bayes factor K can then be interpreted with respect to how strongly the null hypothesis is rejected by the data.

The current version of PhyloSift supports application of Bayesian hypothesis tests to a concatenated alignment of the 37 elite gene families or any other single marker gene, and can be applied to phylogenies inferred either from amino acid or codon-aligned DNA sequences.

Community structure comparison: application to human microbiome data

In addition to hypothesis testing for lineages, PhyloSift also provides a platform to conduct comparative analysis of microbial community structure directly from metagenomic data. To understand how community structure analysis with PhyloSift compares to similar analysis based on 16S rRNA amplicon sequencing we study a recently published human microbiome dataset where samples were sequenced both by a 16S amplicon and a shotgun metagenome approach [67]. In that study, fecal material was collected from infants and adults at diverse geographical locations and subjected to sequencing. Over 600 samples were sequenced using the 16S amplicon protocol. Of those 106 were also subjected to metagenomic shotgun sequencing using 454 pyrosequencing chemistry. Here we apply PhyloSift to the 106 metagenomic samples and conduct a community structure comparison among the samples, and replicate the Yatsunenko *et al.* QIIME analyses on this subset of data.

All QIIME analyses were carried out using release 1.5.0 of the QIIME software toolkit, using the workflow and parameters reported by Yatsunenko *et al.* The Greengenes reference database (collapsed at 97% identity) was used to carry out a closed-reference OTU picking protocol at 97% sequence identity with uclust. All reads which matched database sequences at this level were retained for downstream processing, while non-matching sequences were excluded from further analyses. Parameters for the pick_otus.py script were as follows: `-max_accepts 1 -max_rejects 8 -stepwords 8 -word_length 8`. Taxonomic assignments for OTUs were given by the Greengenes database. Rarefaction and PCoA analyses were carried out using the `alpha_diversity.py` and `beta_diversity_through_plots.py` workflows. A full list of these QIIME commands and output files have been publicly deposited in figshare (DOI: 10.6084/m9.figshare.650869).

PhyloSift processed each of the 106 samples, requiring an average of 2.5 hours per sample on a single 2.27GHz Intel Xeon E5520 core (circa 2009 model). The majority of CPU time is spent in phylogenetic placement of reads. These samples have 154,485 non-human sequence reads on average, for an average of 52 Mbp of sequence data per sample.

We then conducted Edge Principal Components Analysis (PCA) using the reads placed onto the phylogeny of elite gene families. Edge PCA identifies the combination of phylogenetic lineages that explain the greatest extent of variation in the microbial communities in each sample. The resulting PCA plot is shown in Figure 2, with each sample colored according to the age of the human host at the time of sampling. The PCA reveals a strong association between age and microbial community structure. This relationship was also identified by Yatsunenko *et al* using 16S rRNA analysis on a set of >600 samples which included the 106 studied here. In order to quantify the degree of similarity between the PhyloSift Edge PCA and QIIME PCoA results, we calculated Procrustes distances among each pair of analyses, the results are given in Table 1. In general we find that QIIME's PCoA analysis of metagenomic 16S reads produces results that are very different to all other methods, whereas results produced by QIIME

	QIIME 16S Meta	PhyloSift 16S Meta	PhyloSift Elite Meta
QIIME 16S Amp	0.5134279	0.3873677	0.3762175
QIIME 16S Meta	-	0.5376786	0.6351224
PhyloSift 16S Meta	-	-	0.2450837

Table 1. Procrustes distances between microbial community analysis methods. Analysis of 16S amplicon sequences with QIIME (QIIME 16S Amp) produces results more similar to PhyloSift analyzing either 16S or elite protein sequences from metagenomic data than to QIIME analysis of 16S sequences from metagenomic data. PhyloSift results for 16S and elite proteins are more similar to each other than to either QIIME method, possibly due to differences between Edge PCA and the QIIME-generated PCoA on UniFrac distances.

PCoA analysis of 16S amplicon data are more similar to results produced by PhyloSift on metagenomic data.

The nature of edge PCA lends itself to an intuitive inspection of the phylogenetic lineages explaining the difference in community structures. PhyloSift, by using pplacer's guppy program and the Archaeopteryx tree viewer, can produce a visualization of the lineages most strongly associated with each principal component. Figure 3 shows this visualization for the edge PCA analysis of 106 fecal metagenome communities. In that figure, lineages are thickened proportionally to their contribution to the principal component, and are colored according to whether they increase (red) or decrease (turquoise) in abundance along the principal component axis. As we can see from Figure 3 left, the first principal component is defined by an increase in Ruminococcaceae, Clostridiales, and Bacteroides, with a decrease in Bifidobacteria. The association with age suggests that as communities develop in aging children, the Bifidobacteria become less abundant and members of those other lineages grow in abundance. The analysis of Yatsunenکو *et al* on 16S rRNA data also identified age-associated increases in Ruminococcaceae and Bacteroides and a decrease in Bifidobacteria.

Whereas the first principal component agrees strongly with the analysis reported by Yatsunenکو *et al*, the second principal component appears to identify a previously unreported aspect of variation in these samples. Extreme samples on the 2nd principal component (PC2) are very young infants whose fecal microbiota appear to be dominated not by Bifidobacteria, but instead by members of the genus Enterobacter and family Lactobacillales (see Figure 3, right). One possible explanation for this observation may be an association with breast-feeding status of the infants. However, inspection of publicly available metadata did not reveal any clear association of PC2 with breastfeeding status or other recorded metadata. Another possible explanation is mode of birth, vaginal or caesarian, however no information on mode of birth is available for this dataset (Jeffrey Gordon, personal communication). We note that members of the Lactobacillales are abundant in the human vaginal tract, suggesting that newborns high on the 2nd principal component axis may be vaginally delivered if the two groups of newborns do indeed reflect differences in mode of delivery. Interestingly, the dimensions of community structure variation identified in the current set of 106 samples differ from those identified by Yatsunenکو *et al* in the larger set of 600 samples for which amplicon data are available. Geography and age were associated with most variation in their analysis of >600 samples, and the 106 metagenome samples are primarily from infants and do not equally represent that variation. It seems that age-related variation in the microbiome dominates the 106 metagenome samples.

We also investigated the diversity of microbes in the fecal samples. Classic measures of species diversity such as alpha and beta diversity have been applied to microbial communities by collapsing sequences to operational taxonomic units (OTUs). More recently, phylogenetic diversity (PD) [20] has been applied to metagenomic data, yielding a diversity metric that does not require defining OTUs [27]. In the present work we compute phylogenetic diversity on the placed reads, using the attachment points of reads to the reference tree as the basis for the diversity calculation. Figure 4 shows the phylogenetic

diversity present in the fecal samples as a function of age. We observe a general trend where phylogenetic diversity grows quickly with age, presumably due to colonization of the infant gut, then continues to grow slowly throughout adult life. There is a significant log-linear relationship of phylogenetic diversity with age (Pearson's product-moment correlation, $p < 10^{-15}$). We also plot a variant of the PD metric called balance-weighted phylogenetic diversity [38], where diversity contributed by each lineage is weighted by its relative abundance. Balance-weighted PD exhibits a similar growth in early life, but values for individual samples shift relative to population median values. Notably, balance-weighted PD declines in old age, suggesting that a smaller number of divergent lineages may come to dominate the adult human gut. The maximum balance-weighted PD value observed among any sample in the dataset was at the 7th month of life. When samples from before and after the 7th month of life are tested separately, balance-weighted PD exhibits significant age-associated growth before the 7th month ($p = 0.009$, Spearman's rank correlation) and age-associated decline after the 7th month ($p < 10^{-5}$, Spearman's rank correlation). It is not clear what drives the reduction in balance-weighted PD after the 7th month of life, though we note that solid food is commonly introduced to the infant's diet around this time.

PhyloSift provides a means to visualize the relative abundance of taxonomic groups present in a sample. Figure 5 shows two such plots for samples from a 1 month old breastfeeding infant and a 45 year old mother from the Yatsunenkeno *et al* data [67].

Computational efficiency

When processing large metagenomic datasets, computational efficiency and resources can become a logistical challenge. For Illumina data, PhyloSift can process sequence reads on a single CPU core at least as quickly as they can be generated by current instruments. Figure 6 gives memory and running time requirements for some test Illumina datasets. The majority of PhyloSift's running time is spent in phylogenetic read placement (data not shown). Most stages of the workflow implemented by PhyloSift are amenable to both fine and coarse-grain parallelism, thus parallel implementations of the workflow could be created should future data volumes demand it. Finally, the peak memory usage recorded during each run remains roughly constant at 6-9 GB across all data set sizes. As such, PhyloSift is memory-efficient enough to process metagenomic datasets on modern laptop hardware, wherein configurations with 8 GB RAM are readily available.

Discussion

We have presented a new approach for phylogenetic analysis of genomes and uncultured microbial communities. The software implementation of our method, called PhyloSift, also provides a platform for comparison of community structure among many samples. Phylogenetic analysis (placement of short sequences onto reference phylogenies) offers a number of conceptual advantages over OTU-based or taxonomic analysis (interpreting sequence data on the basis of hierarchal classification information) for metagenomic data. Without applying phylogenetic analysis, taxonomic analysis can produce results that are difficult to interpret, particularly when an unknown environmental sequence contains many high scoring hits to reference database sequences as is common in BLAST-based approaches. Alternatively, taxonomic information can be misleading for sequences from species lacking close relatives in public sequence databases; these sequences may recover no match at all, or be assigned taxonomic annotations which do not accurately reflect phylogenetic relationships (e.g. the closest match is still a distant relative, as reflected by low BLAST scores) [16]. Phylogenetic analysis avoids both of these problems, relying instead on evolutionary models to accurately place unknown sequences within a known topology. In many cases, phylogenies will also offer a higher resolution representation of genetic ancestry than taxonomies. For these reasons, we focus on types of phylogenetic analysis enabled by PhyloSift and forgo a discussion of previous taxonomy-based metagenome analysis methods.

Phylogenetic analysis of metagenome sequence data could in principle offer several advantages in the area of microbial forensics. First, by studying an uncultured community, some potential pitfalls of culture bias and sample contamination can be avoided entirely. Second, the environmental shotgun sequencing approach can avoid problems related to PCR primer bias, though issues related to DNA extraction bias remain a problem [44] and might be especially relevant for sporulating organisms such as the Bacilli and their relatives. Third, the metagenomic approach can be applied without prior knowledge of which genes to target in the sample, and permits interrogation of both slow-evolving genes such as 16S rRNA and fast evolving genes that might offer greater resolution among closely related organisms. Finally, phylogenetics can be applied to any gene of interest regardless of whether its evolutionary history is concordant with a taxonomic hierarchy.

Here we have introduced a means to statistically test for lineages of interest directly from an uncultured DNA sample. The test calculates a Bayes factor for the two competing hypotheses: zero sequences derive from the target lineage, versus one or more sequences in the sample derive from the target lineage. This method can be applied to any protein-coding or noncoding gene family of interest. Certain gene families will yield more sensitive tests than others, for example the 16S rRNA gene is slow-evolving and can not usually distinguish within-species relationships where some protein-coding genes might have greater resolution. We emphasize that the Bayes factor is not a test of homology – homology tests exist as e-value and related score statistics in aligners such as BLAST, LAST, and HMMER. Given sequences homologous to a gene family, the Bayes factor tests from which lineage they diverged. The limit of detection for this method will depend on how deeply a sample has been sequenced. This value will depend on several factors specific to the sequencing chemistry and currently must be calculated independently by the user.

The 37 elite gene families were selected because they are universally present and almost always in single copy, but there are some exceptions. When partial homologs exist interpretation of the lineage test can become complicated by paralogs or ancient lateral gene transfer events. Thus one must exercise appropriate caution when interpreting the results of the lineage test. It is a test of whether the sample is void of DNA predicted to have derived from a particular lineage in the phylogeny. For applications like medical diagnostics a more elaborate Bayesian hypothesis test might be appropriate. Such a test might check for a collection of genes that are diagnostic of the organism rather than seeking a single gene, based on prior knowledge that most of the 37 genes are present in most lineages. Such an approach would be less sensitive to sporadic lateral gene transfer events in any single gene family and represents a direction for future work.

Although we do not provide examples, it is possible to test the hypothesis that two microbial communities have equal composition using the phylogenetic Kantorovich-Rubenstein distance [19]. In a bioforensics context this approach could be applied to test whether two uncultured communities of interest “match” each other. The implementation of the method employs an efficient approximation to calculate p -values for the null hypothesis of equal community composition and has been described elsewhere [19]. This test can be applied directly to any individual gene family processed by PhyloSift or to the concatenated alignment of elite families at either the amino acid or DNA sequence level. One limitation of this test is that it does not currently provide a means to account for variability in apparent community structure introduced by normal sample handling procedures. Future work might develop tests that employ many technical replicates of samples to account for such variation in the hypothesis test.

PhyloSift can also be applied to explore the variation in community structure present in a collection metagenomic samples. In recent years it has become standard practice to explore microbial community structure variation using amplicon sequencing of highly conserved genes such as 16S rRNA, 18S rRNA and ITS regions followed by analysis with a pipeline such as QIIME [11], VAMPS (<http://vamps.mbl.edu>), or mothur [50]. Analysis of community structure using metagenome sequence has some potential advantages, such as avoiding issues related to PCR primer bias and distinguishing between erroneous PCR chimeras and sequences representing the “rare biosphere” [5]. However, there are also shortcomings, such as the relatively sparse phylogenetic diversity of available reference genomes relative to amplicon databases.

The reference-based approach taken by PhyloSift will suffer this database resolution limitation when processing metagenomic data, although not when processing amplicon data. Efforts to increase the phylogenetic diversity of available genome sequences are ongoing [63, 48, 53]

Advances in the preparation of high throughput samples will make comparative metagenomics more tractable. The analysis we describe of human fecal microbial communities was possible with a median of only 50 Mbp sequence data per sample. Current Illumina HiSeq 2000 instruments generate up to 40 Gbp per lane, suggesting that up to 800 samples could be processed in a single Illumina lane and yield similar findings. Based on current Illumina sequencing service provider costs this suggests large-scale gut metagenome surveys could be conducted for as little as to \$2.50 to \$5 per sample in sequencing costs. Library preparation would dominate the overall cost of such studies, as current kits from Illumina require about \$37 per sample.

Although we focus on phylogenetic analysis in this work, PhyloSift also provides a basic mechanism to attach taxonomic labels to branches of the phylogenetic trees. Our approach for taxonomic labeling of the phylogeny does not enforce a strict 1:1 mapping between taxonomic labels and branches in the phylogeny. Rather, each branch in the phylogeny is labeled with the entire set of most topologically consistent taxonomic labels. In cases where gene trees may be discordant with the taxonomic tree, this approach allows PhyloSift to represent some of this ambiguity in its results. A systematic study investigating the relationship between rates and patterns of LGT and the effectiveness of our approach for taxonomic labeling remains as future work, as does extension of the taxonomic labeling method to gene families for which duplication and loss is prevalent.

One major limitation of the current approach is that all phylogenetic analysis is conducted independently on each gene. However, genes do not evolve in isolation but rather co-evolve with each other in genomes. Recent studies have demonstrated that large parts of the phylogenetic history in different microbial genes are congruent even though they have undergone lateral gene transfer, duplication, and loss [56, 7]. Large-scale statistical inference of phylogenetic networks (e.g. on > 1000 microbial genomes) that account for duplication, transfer, and loss histories have not yet been described in the literature, however if such a network could be constructed it might provide a means to co-analyze all genes and yield a corresponding increase in sensitivity and power for statistical tests.

Availability

Software for Linux and Mac OS X, along with source code is freely available from <http://github.com/gjospin/PhyloSift>. Extensive user documentation is available at <http://phylosift.wordpress.com>. The source code has been licensed under the GNU Public License (GPL) v3.0.

References

- [1] Sophie S. Abby, Eric Tannier, Manolo Gouy, and Vincent Daubin. Lateral gene transfer as a support for the tree of life. *Proceedings of the National Academy of Sciences*, 2012. doi:10.1073/pnas.1116871109. URL <http://www.pnas.org/content/early/2012/03/12/1116871109.abstract>.
- [2] Andrew Adey, Hilary Morrison, Asan, Xu Xun, Jacob Kitzman, Emily Turner, Bethany Stackhouse, Alexandra MacKenzie, Nicholas Caruccio, Xiuqing Zhang, and Jay Shendure. Rapid, low-input, low-bias construction of shotgun fragment libraries by high-density in vitro transposition. *Genome Biology*, 11(12):R119, 2010. ISSN 1465-6906. doi:10.1186/gb-2010-11-12-r119. URL <http://genomebiology.com/content/11/12/R119>.

- 613 [3] Stephen F. Altschul, Thomas L. Madden, Alejandro A. Schffer, Jinghui Zhang, Zheng Zhang, Webb
614 Miller, and David J. Lipman. Gapped BLAST and PSI-BLAST: a new generation of protein database
615 search programs. *Nucleic Acids Research*, 25(17):3389–3402, 1997. doi:10.1093/nar/25.17.3389. URL
616 <http://nar.oxfordjournals.org/content/25/17/3389.abstract>.
- 617 [4] Oded Béja, L Aravind, Eugene V Koonin, Marcelino T Suzuki, Andrew Hadd, Linh P Nguyen,
618 Stevan B Jovanovich, Christian M Gates, Robert A Feldman, John L Spudich, et al. Bacterial
619 rhodopsin: evidence for a new type of phototrophy in the sea. *Science*, 289(5486):1902–1906, 2000.
- 620 [5] H.M. Bik, D.L. Porazinska, S. Creer, J.G. Caporaso, R. Knight, and W.K. Thomas. Sequencing our
621 way towards understanding global eukaryotic biodiversity. *Trends in ecology & evolution*, 2012.
- 622 [6] Paul C Blainey. The future is now: single-cell genomics of bacteria and archaea. *FEMS microbiology*
623 *reviews*, 2013.
- 624 [7] Bastien Boussau, Gergely J. Szllsi, Laurent Duret, Manolo Gouy, Eric Tannier, and Vin-
625 cent Daubin. Genome-scale coestimation of species and gene trees. *Genome Research*,
626 2012. doi:10.1101/gr.141978.112. URL [http://genome.cshlp.org/content/early/2012/11/06/](http://genome.cshlp.org/content/early/2012/11/06/gr.141978.112.abstract)
627 [gr.141978.112.abstract](http://genome.cshlp.org/content/early/2012/11/06/gr.141978.112.abstract).
- 628 [8] A. Brady and S. Salzberg. Phymmbl expanded: confidence scores, custom databases, parallelization
629 and more. *Nature methods*, 8(5):367–367, 2011.
- 630 [9] A. Brady and S.L. Salzberg. Phymm and phymmbl: metagenomic phylogenetic classification with
631 interpolated markov models. *Nature methods*, 6(9):673–676, 2009.
- 632 [10] Christiam Camacho, George Coulouris, Vahram Avagyan, Ning Ma, Jason Papadopoulos, Kevin
633 Bealer, and Thomas Madden. BLAST+: architecture and applications. *BMC Bioinformatics*, 10
634 (1):421, 2009. ISSN 1471-2105. doi:10.1186/1471-2105-10-421. URL [http://www.biomedcentral.](http://www.biomedcentral.com/1471-2105/10/421)
635 [com/1471-2105/10/421](http://www.biomedcentral.com/1471-2105/10/421).
- 636 [11] J.G. Caporaso, J. Kuczynski, J. Stombaugh, K. Bittinger, F.D. Bushman, E.K. Costello, N. Fierer,
637 A.G. Pena, J.K. Goodrich, J.I. Gordon, et al. Qiime allows analysis of high-throughput community
638 sequencing data. *Nature methods*, 7(5):335–336, 2010.
- 639 [12] Kevin Chen and Lior Pachter. Bioinformatics for whole-genome shotgun sequencing of microbial
640 communities. *PLoS computational biology*, 1(2):e24, 2005.
- 641 [13] Naryttza Diaz, Lutz Krause, Alexander Goesmann, Karsten Niehaus, and Tim Nattkemper. TACO-
642 - Taxonomic classification of environmental genomic fragments using a kernelized nearest neighbor
643 approach. *BMC Bioinformatics*, 10(1):56, 2009. ISSN 1471-2105. doi:10.1186/1471-2105-10-56. URL
644 <http://www.biomedcentral.com/1471-2105/10/56>.
- 645 [14] G.J. Dick, A.F. Andersson, B.J. Baker, S.L. Simmons, B.C. Thomas, A.P. Yelton, J.F. Banfield,
646 et al. Community-wide analysis of microbial genome sequence signatures. *Genome Biol*, 10(8):R85,
647 2009.
- 648 [15] Sean R. Eddy. Accelerated Profile HMM Searches. *PLoS Comput Biol*, 7(10):e1002195, 10 2011.
649 doi:10.1371/journal.pcbi.1002195. URL <http://dx.doi.org/10.1371/journal.pcbi.1002195>.
- 650 [16] Jonathan A Eisen. Phylogenomics: improving functional predictions for uncharacterized genes by
651 evolutionary analysis. *Genome research*, 8(3):163–167, 1998.
- 652 [17] Jonathan A Eisen. Environmental shotgun sequencing: its potential and challenges for studying the
653 hidden world of microbes. *PLoS biology*, 5(3):e82, 2007.

- 654 [18] Jonathan A Eisen. Phylogenetic and phylogenomic approaches to analysis of microbial commu-
655 nities. In *The Social Biology of Microbial Communities A Report from the National Academy
656 of Sciences Forum on Microbial Threats*, pages 180–212. National Academy of Sciences, 2012.
657 doi:10.6084/m9.figshare.841773.
- 658 [19] Steven N. Evans and Frederick A. Matsen. The phylogenetic Kantorovich-Rubinstein metric for
659 environmental sequence samples. *Journal of the Royal Statistical Society: Series B (Statistical
660 Methodology)*, 74(3):569–592, 2012. ISSN 1467-9868. doi:10.1111/j.1467-9868.2011.01018.x. URL
661 <http://dx.doi.org/10.1111/j.1467-9868.2011.01018.x>.
- 662 [20] Daniel P. Faith. Conservation evaluation and phylogenetic diversity. *Biological Conservation*, 61(1):1
663 – 10, 1992. ISSN 0006-3207. doi:10.1016/0006-3207(92)91201-3. URL [http://www.sciencedirect.
664 com/science/article/pii/0006320792912013](http://www.sciencedirect.com/science/article/pii/0006320792912013).
- 665 [21] T.S. Ghosh, M.H. Mohammed, D. Komanduri, and S.S. Mande. Provide: A software tool for accurate
666 estimation of viral diversity in metagenomic samples. *Bioinformatics*, 6(2):91–4, 2011.
- 667 [22] Fabio Gori, Gianluigi Folino, Mike S. M. Jetten, and Elena Marchiori. MTR: taxonomic annotation of
668 short metagenomic reads using clustering at multiple taxonomic ranks. *Bioinformatics*, 27(2):196–
669 203, 2011. doi:10.1093/bioinformatics/btq649. URL [http://bioinformatics.oxfordjournals.
670 org/content/27/2/196.abstract](http://bioinformatics.oxfordjournals.org/content/27/2/196.abstract).
- 671 [23] Philip Hugenholtz, Brett M Goebel, and Norman R Pace. Impact of culture-independent studies on
672 the emerging phylogenetic view of bacterial diversity. *Journal of bacteriology*, 180(18):4765–4774,
673 1998.
- 674 [24] Daniel H. Huson, Alexander F. Auch, Ji Qi, and Stephan C. Schuster. MEGAN analysis of
675 metagenomic data. *Genome Research*, 17(3):377–386, 2007. doi:10.1101/gr.5969107. URL [http:
676 //genome.cshlp.org/content/17/3/377.abstract](http://genome.cshlp.org/content/17/3/377.abstract).
- 677 [25] Frederick A. Matsen IV and Steven N. Evans. Edge principal components and squash clustering:
678 using the special structure of phylogenetic placement data for sample comparison. *PLOS ONE*, 8:
679 e56859, 3 2013. doi:10.1371/journal.pone.0056859.
- 680 [26] Keith A. Jolley, Carly M. Bliss, Julia S. Bennett, Holly B. Bratcher, Carina Brehony, Frances M.
681 Colles, Helen Wimalarathna, Odile B. Harrison, Samuel K. Sheppard, Alison J. Cody, and Martin
682 C. J. Maiden. Ribosomal multilocus sequence typing: universal characterization of bacteria from
683 domain to strain. *Microbiology*, 158(Pt 4):1005–1015, 2012. doi:10.1099/mic.0.055459-0. URL
684 http://mic.sgmjournals.org/content/158/Pt_4/1005.abstract.
- 685 [27] Steven W. Kembel, Jonathan A. Eisen, Katherine S. Pollard, and Jessica L. Green. The Phylogenetic
686 Diversity of Metagenomes. *PLoS ONE*, 6(8):e23214, 08 2011. doi:10.1371/journal.pone.0023214.
687 URL <http://dx.doi.org/10.1371/journal.pone.0023214>.
- 688 [28] Szymon M. Kiebas, Raymond Wan, Kengo Sato, Paul Horton, and Martin C. Frith. Adaptive
689 seeds tame genomic sequence comparison. *Genome Research*, 2011. doi:10.1101/gr.113985.110. URL
690 <http://genome.cshlp.org/content/early/2011/02/04/gr.113985.110.abstract>.
- 691 [29] Victor Kunin, Alex Copeland, Alla Lapidus, Konstantinos Mavromatis, and Philip Hugenholtz. A
692 bioinformatician’s guide to metagenomics. *Microbiology and Molecular Biology Reviews*, 72(4):557–
693 578, 2008.
- 694 [30] Jenna Morgan Lang, Aaron E Darling, and Jonathan A Eisen. Phylogeny of bacterial and archaeal
695 genomes using conserved genes: Supertrees and supermatrices. *PloS one*, 8(4):e62510, 2013.

- 696 [31] Ben Langmead, Cole Trapnell, Mihai Pop, and Steven Salzberg. Ultrafast and memory-efficient
697 alignment of short DNA sequences to the human genome. *Genome Biology*, 10(3):R25, 2009. ISSN
698 1465-6906. doi:10.1186/gb-2009-10-3-r25. URL <http://genomebiology.com/2009/10/3/R25>.
- 699 [32] Roger S Lasken. Genomic sequencing of uncultured microorganisms from single cells. *Nature Reviews*
700 *Microbiology*, 10(9):631–640, 2012.
- 701 [33] B. Liu, T. Gibbons, M. Ghodsi, and M. Pop. Metaphyer: Taxonomic profiling for metagenomic
702 sequences. In *Bioinformatics and Biomedicine (BIBM), 2010 IEEE International Conference on*,
703 pages 95–100. IEEE, 2010.
- 704 [34] Ari Löytynoja, Albert J. Vilella, and Nick Goldman. Accurate extension of multiple se-
705 quence alignments using a phylogeny-aware graph algorithm. *Bioinformatics*, 28(13):1684–1691,
706 2012. doi:10.1093/bioinformatics/bts198. URL [http://bioinformatics.oxfordjournals.org/](http://bioinformatics.oxfordjournals.org/content/28/13/1684.abstract)
707 [content/28/13/1684.abstract](http://bioinformatics.oxfordjournals.org/content/28/13/1684.abstract).
- 708 [35] C. Lozupone and R. Knight. Unifrac: a new phylogenetic method for comparing microbial commu-
709 nities. *Applied and environmental microbiology*, 71(12):8228–8235, 2005.
- 710 [36] F.A. Matsen, N.G. Hoffman, A. Gallagher, and A. Stamatakis. A format for phylogenetic placements.
711 *PLOS ONE*, 7(2):e31009, 2012. doi:10.1371/journal.pone.0031009.
- 712 [37] Frederick Matsen, Robin Kodner, and E Virginia Armbrust. pplacer: linear time maximum-likelihood
713 and Bayesian phylogenetic placement of sequences onto a fixed reference tree. *BMC Bioinformatics*,
714 11(1):538, 2010. ISSN 1471-2105. doi:10.1186/1471-2105-11-538. URL [http://www.biomedcentral.](http://www.biomedcentral.com/1471-2105/11/538)
715 [com/1471-2105/11/538](http://www.biomedcentral.com/1471-2105/11/538).
- 716 [38] Connor O McCoy, IV Matsen, and A Frederick. Abundance-weighted phylogenetic diversity mea-
717 sures distinguish microbial community states and are robust to sampling depth. *arXiv preprint*
718 *arXiv:1305.0306*, 2013.
- 719 [39] A.C. McHardy, H.G. Martín, A. Tsirigos, P. Hugenholtz, and I. Rigoutsos. Accurate phylogenetic
720 classification of variable-length dna fragments. *Nature methods*, 4(1):63–72, 2006.
- 721 [40] F Meyer, D Paarmann, M D’souza, R Olson, EM Glass, M Kubal, T Paczian, A Rodriguez,
722 R Stevens, A Wilke, et al. The metagenomics rast server—a public resource for the automatic
723 phylogenetic and functional analysis of metagenomes. *BMC bioinformatics*, 9(1):386, 2008.
- 724 [41] Christopher Miller, Brett Baker, Brian Thomas, Steven Singer, and Jillian Banfield. EMIRGE:
725 reconstruction of full-length ribosomal genes from microbial community short read sequencing data.
726 *Genome Biology*, 12(5):R44, 2011. ISSN 1465-6906. doi:10.1186/gb-2011-12-5-r44. URL [http:](http://genomebiology.com/2011/12/5/R44)
727 [//genomebiology.com/2011/12/5/R44](http://genomebiology.com/2011/12/5/R44).
- 728 [42] M.H. Mohammed, S. Chadaram, D. Komanduri, T.S. Ghosh, and S.S. Mande. Eu-detect: An
729 algorithm for detecting eukaryotic sequences in metagenomic data sets. *J Biosci*, 36(4):709–17,
730 2011.
- 731 [43] M. Monzoorul Haque, Tarini Shankar Ghosh, Dinakar Komanduri, and Sharmila S. Mande. SOrt-
732 ITEMS: Sequence orthology based approach for improved taxonomic estimation of metagenomic
733 sequences. *Bioinformatics*, 25(14):1722–1730, 2009. doi:10.1093/bioinformatics/btp317. URL [http:](http://bioinformatics.oxfordjournals.org/content/25/14/1722.abstract)
734 [//bioinformatics.oxfordjournals.org/content/25/14/1722.abstract](http://bioinformatics.oxfordjournals.org/content/25/14/1722.abstract).
- 735 [44] Jenna L. Morgan, Aaron E. Darling, and Jonathan A. Eisen. Metagenomic Sequenc-
736 ing of an *In Vitro*-Simulated Microbial Community. *PLoS ONE*, 5(4):e10209, 04 2010.
737 doi:10.1371/journal.pone.0010209. URL <http://dx.doi.org/10.1371%2Fjournal.pone.0010209>.

- 738 [45] Brian Ondov, Nicholas Bergman, and Adam Phillippy. Interactive metagenomic visualization in a
739 Web browser. *BMC Bioinformatics*, 12(1):385, 2011. ISSN 1471-2105. doi:10.1186/1471-2105-12-385.
740 URL <http://www.biomedcentral.com/1471-2105/12/385>.
- 741 [46] Kaustubh R Patil, Peter Haider, Phillip B Pope, Peter J Turnbaugh, Mark Morrison, Tobias Scheffer,
742 and Alice C McHardy. Taxonomic metagenome sequence assignment with structured output models.
743 *Nature Methods*, pages 191–192, 2011.
- 744 [47] Morgan N. Price, Paramvir S. Dehal, and Adam P. Arkin. FastTree 2 Approx-
745 imately Maximum-Likelihood Trees for Large Alignments. *PLoS ONE*, 5(3):e9490, 03 2010.
746 doi:10.1371/journal.pone.0009490. URL <http://dx.doi.org/10.1371%2Fjournal.pone.0009490>.
- 747 [48] Christian Rinke, Patrick Schwientek, Alexander Sczyrba, Natalia N Ivanova, Iain J Anderson, Jan-
748 Fang Cheng, Aaron Darling, Stephanie Malfatti, Brandon K Swan, Esther A Gies, et al. Insights
749 into the phylogeny and coding potential of microbial dark matter. *Nature*, 2013.
- 750 [49] Gail L. Rosen, Erin R. Reichenberger, and Aaron M. Rosenfeld. NBC: the Nave Bayes Classification
751 tool webserver for taxonomic classification of metagenomic reads. *Bioinformatics*, 27(1):127–129,
752 2011. doi:10.1093/bioinformatics/btq619. URL [http://bioinformatics.oxfordjournals.org/
753 content/27/1/127.abstract](http://bioinformatics.oxfordjournals.org/content/27/1/127.abstract).
- 754 [50] Patrick D Schloss, Sarah L Westcott, Thomas Ryabin, Justine R Hall, Martin Hartmann, Emily B
755 Hollister, Ryan A Lesniewski, Brian B Oakley, Donovan H Parks, Courtney J Robinson, et al. Intro-
756 ducing mothur: open-source, platform-independent, community-supported software for describing
757 and comparing microbial communities. *Applied and environmental microbiology*, 75(23):7537–7541,
758 2009.
- 759 [51] N. Segata, L. Waldron, A. Ballarini, V. Narasimhan, O. Jousson, and C. Huttenhower. Metagenomic
760 microbial community profiling using unique clade-specific marker genes. *Nature Methods*, 9(8):811–
761 814, 2012.
- 762 [52] T.J. Sharpton, S.J. Riesenfeld, S.W. Kembel, J. Ladau, J.P. O’Dwyer, J.L. Green, J.A. Eisen, and
763 K.S. Pollard. Phylotu: a high-throughput procedure quantifies microbial community diversity and
764 resolves novel taxa from metagenomic data. *PLoS computational biology*, 7(1):e1001061, 2011.
- 765 [53] Patrick M Shih, Dongying Wu, Amel Latifi, Seth D Axen, David P Fewer, Emmanuel Talla, Alexan-
766 dra Calteau, Fei Cai, Nicole Tandeau de Marsac, Rosmarie Rippka, et al. Improving the coverage
767 of the cyanobacterial phylum using diversity-driven genome sequencing. *Proceedings of the National
768 Academy of Sciences*, 110(3):1053–1058, 2013.
- 769 [54] Manuel Stark, Simon Berger, Alexandros Stamatakis, and Christian von Mering. MLTreeMap - accu-
770 rate Maximum Likelihood placement of environmental DNA sequences into taxonomic and functional
771 reference phylogenies. *BMC Genomics*, 11(1):461, 2010. ISSN 1471-2164. doi:10.1186/1471-2164-
772 11-461. URL <http://www.biomedcentral.com/1471-2164/11/461>.
- 773 [55] Shinichi Sunagawa, Daniel R Mende, Georg Zeller, Fernando Izquierdo-Carrasco, Simon A Berger,
774 Jens Roat Kultima, Luis Pedro Coelho, Manimozhiyan Arumugam, Julien Tap, Henrik Bjørn Nielsen,
775 et al. Metagenomic species profiling using universal phylogenetic marker genes. *Nature Methods*,
776 2013.
- 777 [56] Gergely J. Szllsi, Bastien Boussau, Sophie S. Abby, Eric Tannier, and Vincent Daubin. Phylogenetic
778 modeling of lateral gene transfer reconstructs the pattern and relative timing of speciations. *Proceed-
779 ings of the National Academy of Sciences*, 109(43):17513–17518, 2012. doi:10.1073/pnas.1202997109.
780 URL <http://www.pnas.org/content/109/43/17513.abstract>.

- 781 [57] Torsten Thomas, Jack Gilbert, Folker Meyer, et al. Metagenomics-a guide from sampling to data
782 analysis. *Microb Inform Exp*, 2(3), 2012.
- 783 [58] Susannah Green Tringe, Christian Von Mering, Arthur Kobayashi, Asaf A Salamov, Kevin Chen,
784 Hwai W Chang, Mircea Podar, Jay M Short, Eric J Mathur, John C Detter, et al. Comparative
785 metagenomics of microbial communities. *Science*, 308(5721):554–557, 2005.
- 786 [59] Gene W Tyson, Jarrod Chapman, Philip Hugenholtz, Eric E Allen, Rachna J Ram, Paul M Richard-
787 son, Victor V Solovyev, Edward M Rubin, Daniel S Rokhsar, and Jillian F Banfield. Community
788 structure and metabolism through reconstruction of microbial genomes from the environment. *Nature*,
789 428(6978):37–43, 2004.
- 790 [60] J Craig Venter, Karin Remington, John F Heidelberg, Aaron L Halpern, Doug Rusch, Jonathan A
791 Eisen, Dongying Wu, Ian Paulsen, Karen E Nelson, William Nelson, et al. Environmental genome
792 shotgun sequencing of the sargasso sea. *science*, 304(5667):66–74, 2004.
- 793 [61] Q. Wang, G.M. Garrity, J.M. Tiedje, and J.R. Cole. Naive bayesian classifier for rapid assignment
794 of rna sequences into the new bacterial taxonomy. *Applied and environmental microbiology*, 73(16):
795 5261–5267, 2007.
- 796 [62] Tanja Woyke, Damon Tighe, Konstantinos Mavromatis, Alicia Clum, Alex Copeland, Wendy
797 Schackwitz, Alla Lapidus, Dongying Wu, John P. McCutcheon, Bradon R. McDonald, Nancy A.
798 Moran, James Bristow, and Jan-Fang Cheng. One Bacterial Cell, One Complete Genome. *PLoS*
799 *ONE*, 5(4):e10314, 04 2010. doi:10.1371/journal.pone.0010314. URL [http://dx.doi.org/10.1371%](http://dx.doi.org/10.1371%2Fjournal.pone.0010314)
800 [2Fjournal.pone.0010314](http://dx.doi.org/10.1371%2Fjournal.pone.0010314).
- 801 [63] Dongying Wu, Philip Hugenholtz, Konstantinos Mavromatis, Rdiger Pukall, Eileen Dalin, Natalia N
802 Ivanova, Victor Kunin, Lynne Goodwin, Martin Wu, Brian J Tindall, and et al. A phylogeny-
803 driven genomic encyclopaedia of bacteria and archaea. *Nature*, 462(7276):1056–1060, 2009. URL
804 <http://www.ncbi.nlm.nih.gov/pubmed/20033048>.
- 805 [64] Dongying Wu, Guillaume Jospin, and Jonathan A. Eisen. Systematic identification of gene families
806 for use as markers for phylogenetic and phylogeny-driven ecological studies of bacteria and archaea
807 and their major subgroups. *PLoS ONE*, 8(10):e77033, 10 2013. doi:10.1371/journal.pone.0077033.
808 URL <http://dx.doi.org/10.1371%2Fjournal.pone.0077033>.
- 809 [65] M. Wu and A.J. Scott. Phylogenomic analysis of bacterial and archaeal sequences with amphora2.
810 *Bioinformatics*, 28(7):1033–1034, 2012.
- 811 [66] Martin Wu and Jonathan Eisen. A simple, fast, and accurate method of phylogenomic inference.
812 *Genome Biology*, 9(10):R151, 2008. ISSN 1465-6906. doi:10.1186/gb-2008-9-10-r151. URL [http:](http://genomebiology.com/2008/9/10/R151)
813 [//genomebiology.com/2008/9/10/R151](http://genomebiology.com/2008/9/10/R151).
- 814 [67] Tanya Yatsunenko, Federico E. Rey, Mark J. Manary, Indi Trehan, Monica Contreras Maria Glo-
815 ria Dominguez-Bello, Magda Magris, Glida Hidalgo, Robert N. Baldassano, Andrey P. Anokhin,
816 Andrew C. Heath, Barbara Warner, Jens Reeder, Justin Kuczynski, J. Gregory Caporaso, Catherine
817 A. Lozupone, Christian Lauber, Jose Carlos Clemente, Dan Knights, Rob Knight, and Jeffrey I.
818 Gordon. Human gut microbiome viewed across age and geography. *Nature*, 486:222227, 2012.
- 819 [68] Yongan Zhao, Haixu Tang, and Yuzhen Ye. RAPSearch2: a fast and memory-
820 efficient protein similarity search tool for next generation sequencing data. *Bioinformatics*,
821 2011. doi:10.1093/bioinformatics/btr595. URL [http://bioinformatics.oxfordjournals.org/](http://bioinformatics.oxfordjournals.org/content/early/2011/10/28/bioinformatics.btr595.abstract)
822 [content/early/2011/10/28/bioinformatics.btr595.abstract](http://bioinformatics.oxfordjournals.org/content/early/2011/10/28/bioinformatics.btr595.abstract).

823 **Figure Legends**

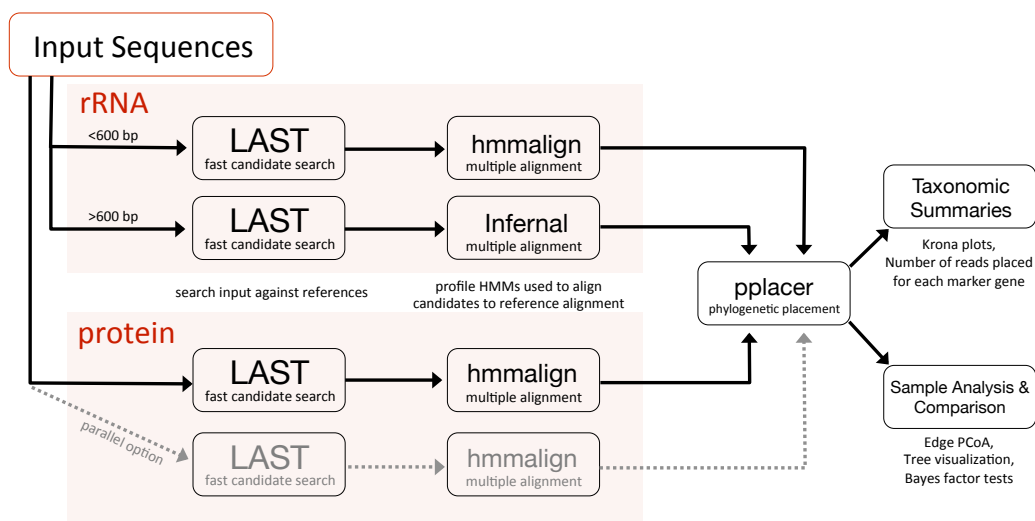


Figure 1. PhyloSift client workflow. This workflow is applied to the user's sequence data. DNA input sequences are processed via both the rRNA and protein parts of the workflow.

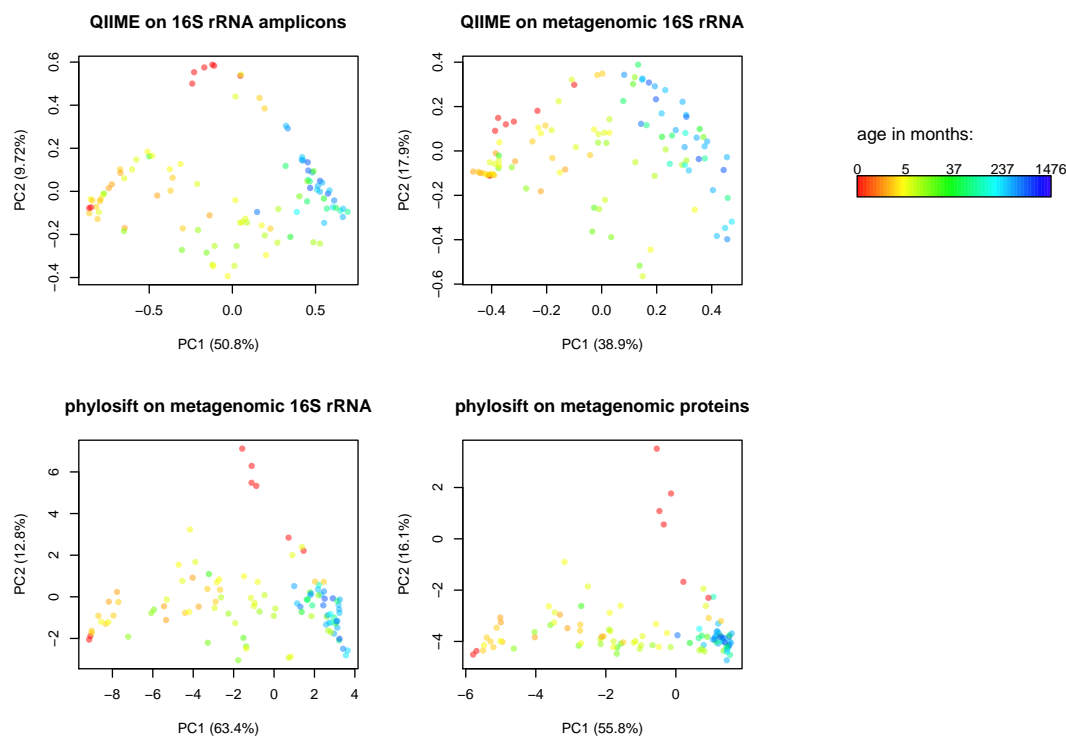


Figure 2. Comparison of QIIME PCA and edge PCA analysis of human fecal samples. Samples from 106 individuals were analyzed by PCA to evaluate trends in community composition with respect to host age. 16S rDNA amplicon data and metagenomic data from the same samples was processed using QIIME and PhyloSift. QIIME analyzed the amplicon data (top left) and 16S rDNA reads extracted from the metagenomic data (top right) using a reference-based OTU picking strategy. PhyloSift analyzed the same metagenomic 16S rDNA reads (bottom left) and reads matching the 37 elite gene families (bottom right). Each PCA approach gives qualitatively similar results, differences as quantified by Procrustes analysis are given in Table 1.

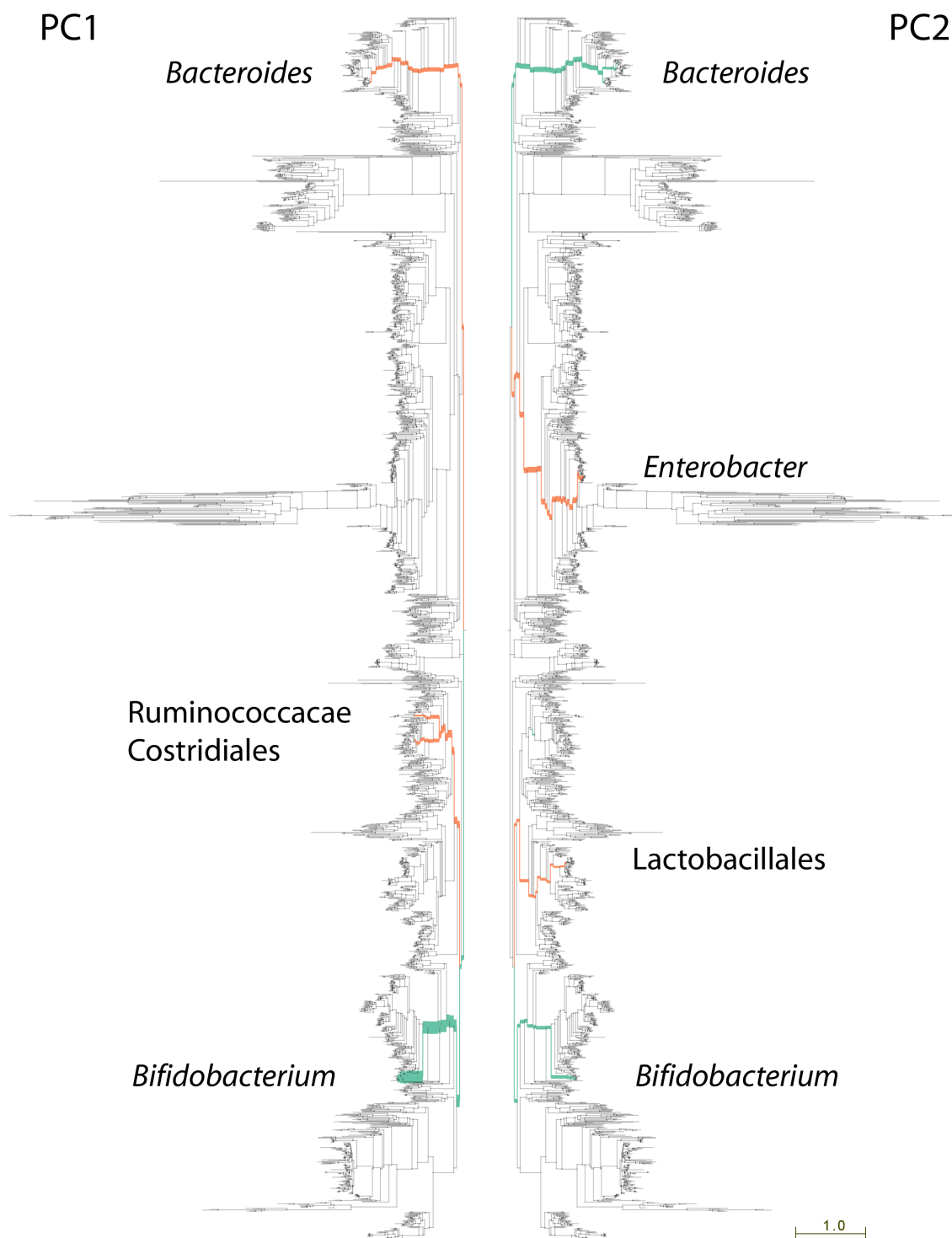


Figure 3. Lineages contributing variation in human fecal sample community structure. 106 metagenomic samples were processed using PhyloSift and their community composition compared using Edge PCA [25]. Lineages that decrease in abundance along the principal component axis are shown in turquoise color, those increasing in abundance are shown in red. Edge width is proportional to the change in abundance. Remaining lineages in the phylogeny of bacteria, archaea, eukarya, and some viruses are shown in light gray. PC1 shown at left, PC2 at right.

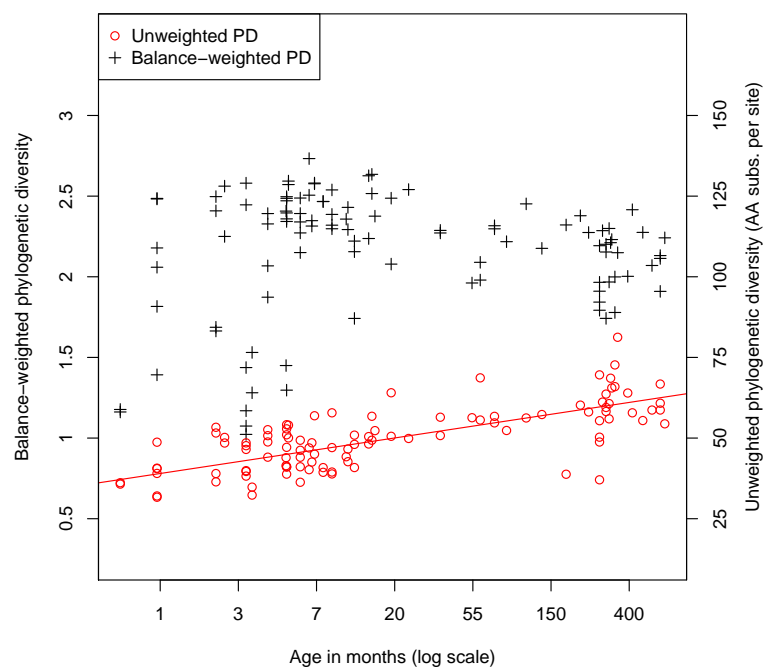


Figure 4. Relationship between fecal community phylogenetic diversity and host age. 106 metagenomic samples were processed using PhyloSift and their phylogenetic diversity analyzed using two metrics. Unweighted phylogenetic diversity (PD) simply measures the total branch length of the reference tree covered by placed reads from a sample. Balance-weighted phylogenetic diversity adjusts these values by the abundance of each lineage in the sample. In unweighted PD, a log-linear relationship between host age and fecal community phylogenetic diversity can be observed. Balance weighted PD, on the other hand, shows rapid growth in early life followed by slow decline after the first year, consistent with a small number of divergent lineages becoming dominant in the fecal ecosystem.

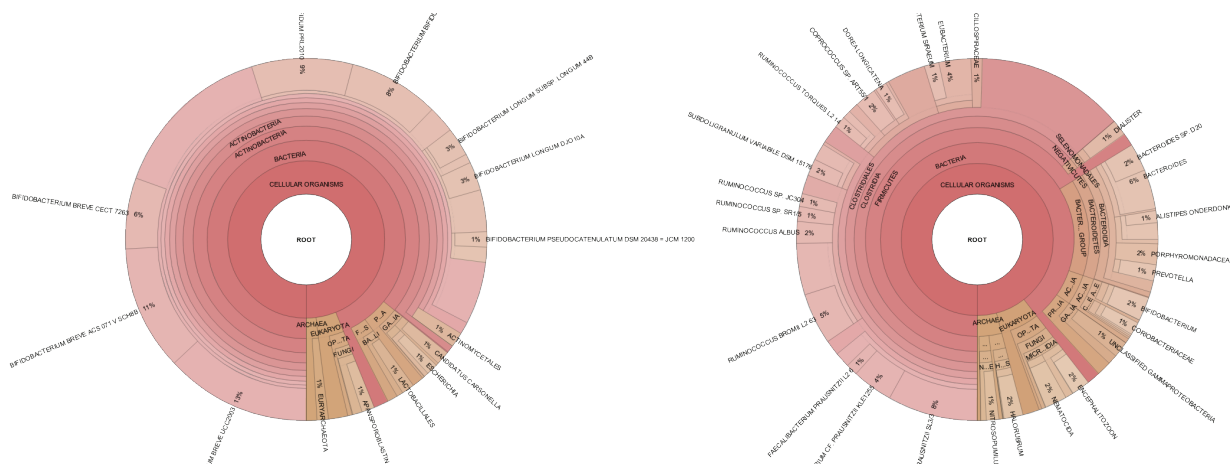


Figure 5. Taxonomic visualization of two human gut samples. Taxonomic plot at left shows an infant, plot at right shows a 45 year old mother. Data analyzed by PhyloSift, visualized by Krona.

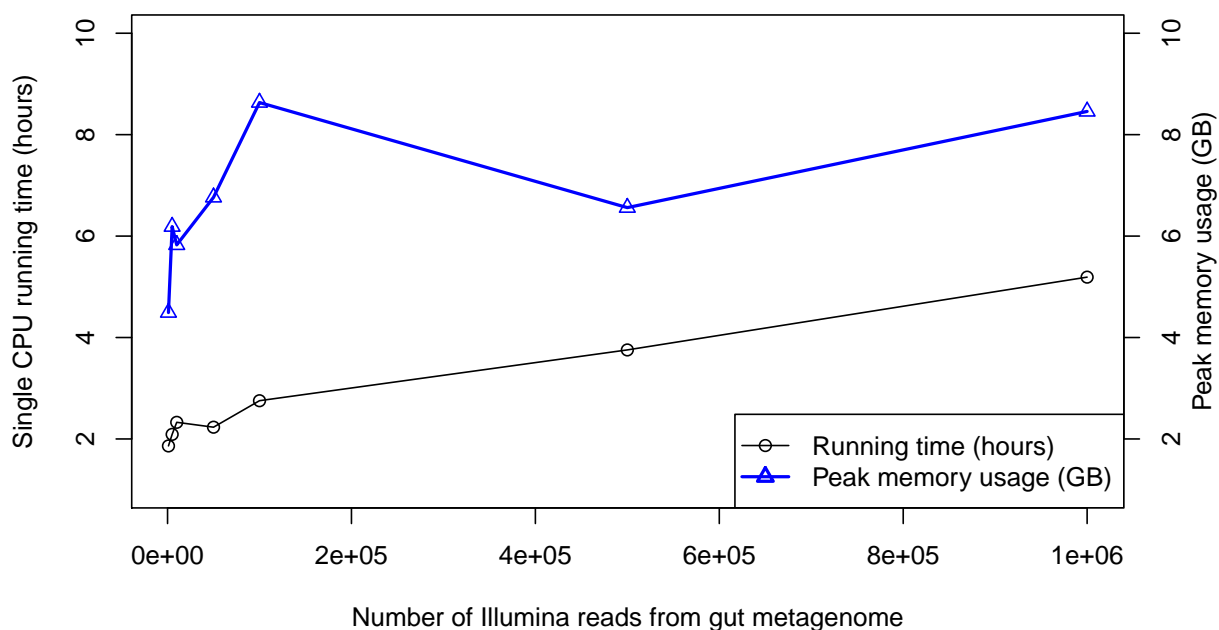


Figure 6. PhyloSift performance and scaling behavior. PhyloSift v1.0 was used to process Illumina sequence data from a human gut microbiome dataset subsampled to varying numbers of reads. The program was run single-threaded on an Intel Xeon E5520 CPU core (circa 2009 model).