

# METAPALETTE: A $k$ -MER PAINTING APPROACH FOR METAGENOMIC TAXONOMIC PROFILING AND QUANTIFICATION OF NOVEL STRAIN VARIATION.

DAVID KOSLICKI<sup>1\*</sup>, DANIEL FALUSH<sup>2</sup>

<sup>1</sup> *Mathematics Department, Oregon State University, Corvallis, OR 97330*

<sup>2</sup> *Institute of Life Sciences, University of Swansea, Singleton Park, Swansea, SA2 8PP UK*

**ABSTRACT.** Metagenomic profiling is challenging in part because of the highly uneven sampling of the tree of life by genome sequencing projects and the limitations imposed by performing phylogenetic inference at fixed taxonomic ranks. We present the algorithm **MetaPalette** which uses long  $k$ -mer sizes ( $k = 30, 50$ ) to fit a  $k$ -mer “palette” of a given sample to the  $k$ -mer palette of reference organisms. By modeling the  $k$ -mer palettes of unknown organisms, the method also gives an indication of the presence, abundance, and evolutionary relatedness of novel organisms present in the sample. The method returns a traditional, fixed-rank taxonomic profile which is shown on independently simulated data to be one of the most accurate to date. Tree figures are also returned that quantify the relatedness of novel organisms to reference sequences and the accuracy of such figures is demonstrated on simulated spike-ins and a metagenomic soil sample.

The software implementing **MetaPalette** is available at:

<https://github.com/dkoslicki/MetaPalette>

Pre-trained databases are included for Archaea, Bacteria, Eukaryota, and viruses.

## 1. INTRODUCTION

Metagenomics is a developing field used to characterize the organismal composition of microbial communities in environmental or clinical samples [45]. A key step in most metagenomic analyses is to identify the organisms in the sample and their relative frequency. A wide variety of different algorithms have been developed for this purpose.

Most approaches, including the one described here, are based on relating sequenced reads to reference organism genome sequences. Conceptually, the aim of these approaches is to place the organisms in the sample on a “tree of life” that has been defined in advance. In practice, the available reference organisms are extremely unevenly scattered through the true tree of life. Many medically important branches, such as enterobacteria, are relatively well sampled with many strains from the same species, while there are entire phyla of unculturable organisms that are unrepresented [3, 25, 32, 40].

A further difficulty, both in theory and in practice, is that a fully resolved tree of life cannot be established, even from complete reference genomes. At the scale of individual species, homologous recombination scrambles variation so that a tree is not necessarily an appropriate representation of organismal relationships, while more distant phylogenetic relationships can be difficult to estimate due to the various technical challenges of reconstructing ancient evolutionary events [38, 39, 41, 43, 51].

---

*Date:* February 16, 2016.

\*Corresponding Author: [david.koslicki@math.oregonstate.edu](mailto:david.koslicki@math.oregonstate.edu).

Based on these practical considerations, an effective metagenomic method should both identify the closest organism or set of organisms in the reference set and also estimate the genetic difference between the closest reference(s) and the organism present in the sample. The method should work both if the closest neighbor is a distant member of the same phyla or if there are multiple strains within the species in question. Fine scale classification is important because the detailed knowledge we have of, for example, *E. coli* shows that organisms from the same species can have entirely different ecology and phenotypic effects on their host [22].

Given these difficulties, a number of different approaches are taken to characterize metagenomic samples. A commonly used approach is first to place individual reads onto a tree constructed for a particular set of genes, and then attempt to sum the phylogenetic information across the reads [16, 30, 36, 47]. Phylogenetic analysis of each read can be computationally challenging for large datasets and individual reads can often only be placed inaccurately. It is challenging to appropriately represent this uncertainty in later stages of the analysis. These approaches also break down if a tree is not a good representation of relationships amongst organisms, e.g. within species. Furthermore, while utilizing specific genes (so-called *marker genes*) can increase computational efficiency, this approach throws away a considerable amount of information from sequences that do not align to the marker genes. As a result of these issues, these methods are typically accurate for genus level or higher classification but not for fine scale classification.

Another approach identifies features that are characteristic of particular organisms, such as the frequency of  $k$ -mers [24, 37, 52]. These features are used either for taxonomic binning of individual reads or in order to compute the overall composition. Depending on the  $k$ -mer size utilized, these methods are either suitable only for higher level phylogenetic analysis (for small  $k$ -mers), or are highly dependent on the training database utilized (for larger  $k$ -mers). In either case, no existent method using this approach can accurately detect and classify organisms highly diverged from ones in the training database, and still struggle with quantifying strain-level variation. Using longer  $k$ -mers allows for higher specificity but using  $k$ -mers that are unique to specific taxa in the reference dataset (as in [37, 52]) ignores a great deal of information about evolutionary relatedness provided by other  $k$ -mers. It also makes the approach highly dependent on the specific composition of the reference dataset. We argue that utilizing all  $k$ -mers in a reference database and multiple  $k$ -mer sizes allows for the modeling of the  $k$ -mer signature of organisms absent from a given training database.

In this manuscript, we present an approach based on defining a “palette” for each reference organism. Specifically, we count the number of  $k$ -mers found in the sample DNA that are present in each reference organism. Our approach thus uses all  $k$ -mers of a particular length in the reference dataset, while discarding the specific information provided by matches of individual  $k$ -mers. This is similar in spirit to the so-called pseudoalignment approach of [44] except here we use  $k$ -mer counts of the entire sample, not of individual reads whose origins may be ambiguous. We model these palettes using a simple linear mixture model which includes both the reference organisms and “hypothetical organisms” at different degrees of genetic relatedness to the reference organisms. The algorithm is called MetaPalette and the outputs of the algorithm are demonstrated in Figure 1.

We first introduce the concept of a *common k-mer matrix* and demonstrate how utilizing multiple  $k$ -mer sizes allows for accurate quantification of evolutionary relatedness. We then develop a mixture modeling procedure that utilizes this information to taxonomically profile a metagenomic sample and indicate the evolutionary relatedness of novel organisms. Evidence on simulated and real data is given that this approach can accurately capture strain-level variation, and we then benchmark this approach against other, commonly utilized metagenomic profiling techniques.

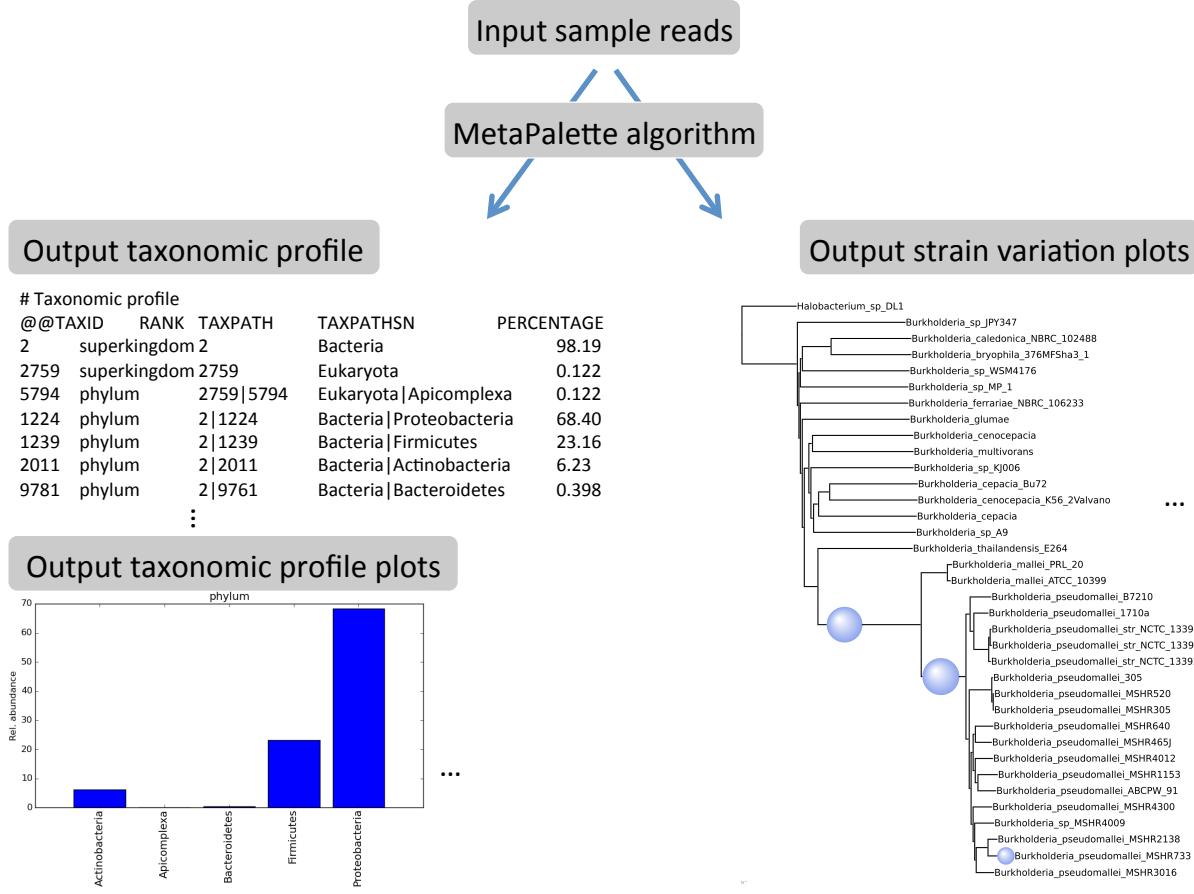


FIGURE 1. Illustration of the MetaPalette algorithm. Along with an output taxonomic profile and bar chart plots at all inferred taxonomic ranks, figures of strain-level variation for each inferred genus and/or species are also output.

## 2. METHODS

**2.1. Common  $k$ -mer Training Matrix.** To quantify the similarity of two genomes, we count (with multiplicity) the fraction of each genome's  $k$ -mers that are in common with the other. Rigorous mathematical definitions of this and other quantities are contained in the Appendix section A. This quantity, denoted  $pckm_k(\cdot, \cdot)$  for *percent common  $k$ -mers*, is similar to the well-known Jaccard index [21] except that, among other differences,  $pckm_k(\cdot, \cdot)$  is not symmetric but does incorporate the counts of  $k$ -mers, not just their occurrence.

When given a set of genomes (i.e. a training database), a pair-wise similarity matrix can be formed:  $A_{i,j}^{(k)} = pckm_k(g_i, g_j)$  for  $g_i$  and  $g_j$  training genomes. The column vector  $pckm_k(\cdot, g_j)$  can be thought of as a palette, representing the particular  $k$ -mer profile of  $g_j$  in relation to other genomes. We call each of these matrices a *common  $k$ -mer matrix*. These matrices reflects the relatedness of the training genomes based on  $k$ -mer similarity. For larger  $k$ -mer sizes, one can clearly extract taxonomic information from these matrices: see Figure 2.

Beyond genera-level variation, strain-level variation can be captured through these common  $k$ -mer matrices. For example, using all the strains of the species *Burkholderia multivorans* accessible via NCBI, we formed a neighbor joining tree using the average of the 30-mer and 50-mer common

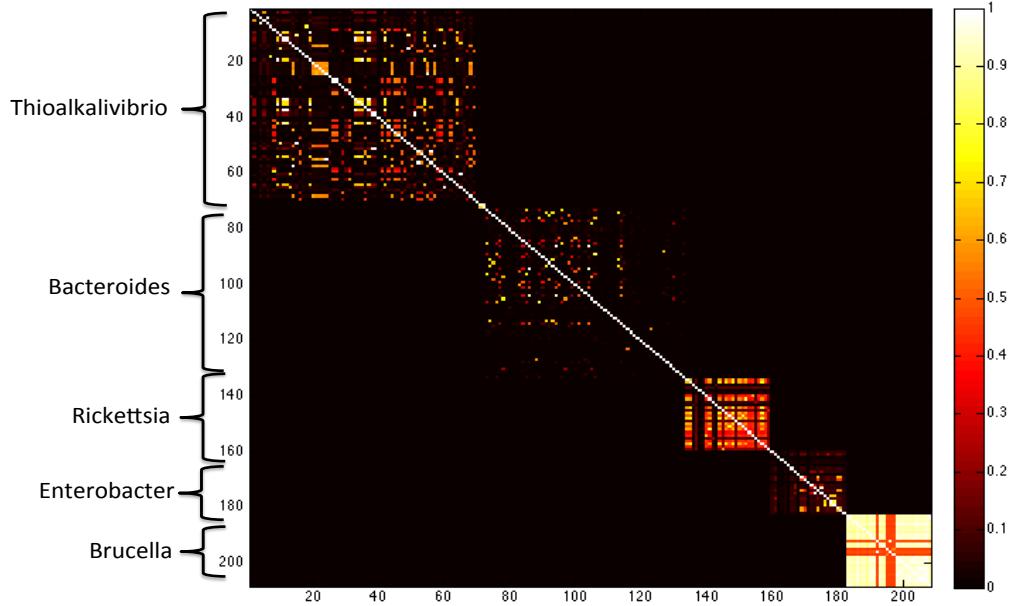


FIGURE 2. Heatmap of the common  $k$ -mer matrix  $A^{(40)}$  for  $k = 40$  using a subset of the NCBI bacterial genome database. Delineations between genera can clearly be seen. In a given genera, differing similarity of species is also visible.

$k$ -mer matrices. This tree, shown in figure 3, demonstrates how the common  $k$ -mer matrices can capture variations amongst these strains.

The entries of  $A^{(k)}$  can be calculated in a computationally efficient manner. We take the approach of forming bloom count filters (using Jellyfish [31]) for each of the training genomes, and then counting the common  $k$ -mers using a simple C++ program based on heap data structures.

**2.2. Modeling Related Organisms.** To model the  $k$ -mer counts for organisms related at varying degrees from the training database, we take advantage of the differing behavior of  $pckm_k(\cdot, \cdot)$  as a function of  $k$  for closely related organisms and distantly related organisms. In particular, the percent of common  $k$ -mers  $pckm_k(\cdot, \cdot)$  decays much slower as a function of  $k$  for closely related organisms than for distantly related ones. This is consistent with the intuition that, for example, two organisms from different phyla will have a similar percent of shared 1-mers, but very few common 50-mers. Conversely, two closely related strains will have both a high percentage of shared 1-mers *and* a high percentage of shared 50-mers. This is demonstrated in Figure 4(a). This property means that using more than one  $k$ -mer length should in principle allow us to distinguish between having an organism that is identical to a training organism at a low frequency and having an organism that is distantly related to all training organisms but present in the sample at a higher frequency.

We focus on two particular  $k$ -mer sizes:  $k = 30$  and  $k = 50$  due to the predictability of  $pckm_k$  for these  $k$ -mer sizes. Indeed, using 6,914 whole bacterial genomes downloaded from a variety of publicly accessible repositories (via RepoPhAn: <https://bitbucket.org/nsegata/repophlan>), we observed that the percent of shared 30-mers can be predicted from the percent of shared 50-mers. See part (b) of Figure 4. A degree 3 polynomial was used (as it resulted in the lowest RMSE

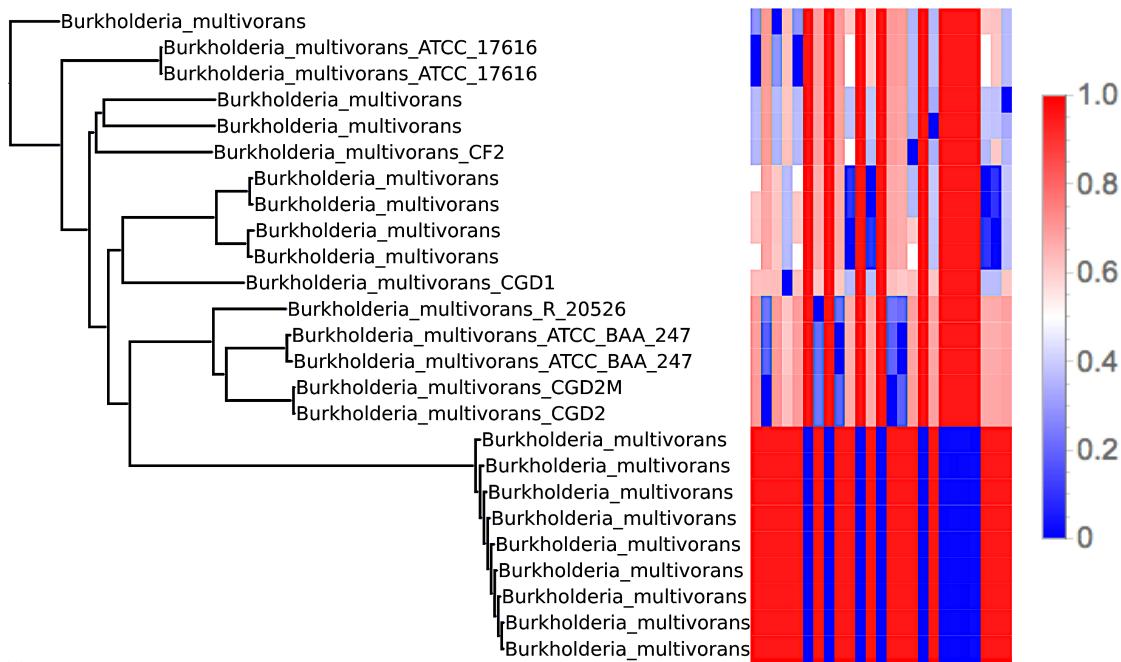


FIGURE 3. Neighbor joining tree for the species *Burkholderia multivorans* based on average of the common 30-mer and 50-mer matrices (shown in heat map to the right) depicting the ability of the common  $k$ -mer matrices to capture strain-level variation.

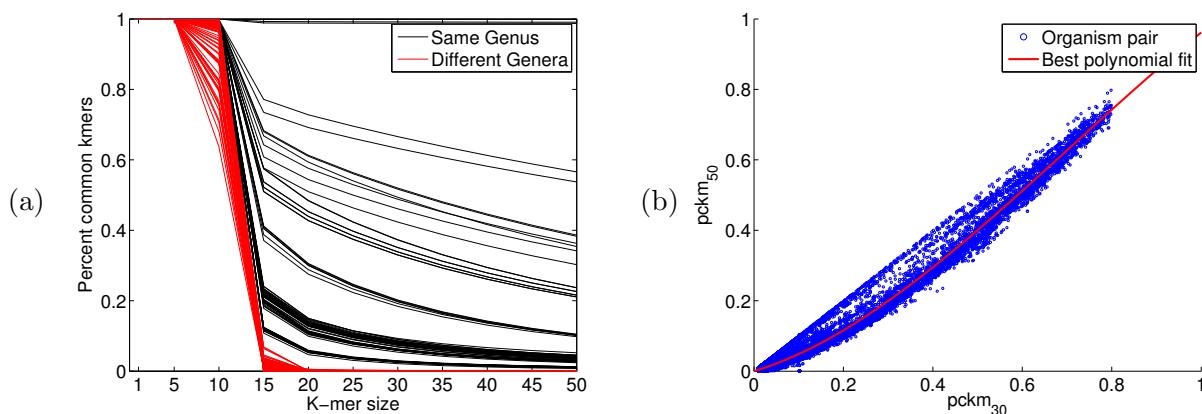


FIGURE 4. (a) Plot of  $k$ -mer similarity  $pckm_k(g_i, g_j)$  as a function of  $k$  for 100 organism pairs of the same genus and 100 of different genera. (b) Scatter plot of the  $6,914^2$  pairs of entries of the common 30-mer and 50-mer matrices shown with the best fit polynomial.

and  $R^2$  values, which did not improve for higher degree polynomials). Namely, we observed that for the polynomial  $p(x) = -.5141x^3 + 1.0932x^2 + 0.3824x$ ,  $pckm_{50}(g_i, g_j) \approx p(pckm_{30}(g_i, g_j))$ .

For  $k$ -mer lengths substantially shorter than 30, the behavior of  $pckm_k$  is more variable, for example because of convergence of sequence composition between distantly related organisms. On the other hand,  $k$ -mers much larger than 50 are increasingly time consuming to compute and are likely to be more sensitive to sequencing error and other technical artifacts.

We can augment the matrices  $A^{(k)}$  with columns that represent hypothetical organisms which are related by different degrees to the reference organism. For a given organism with genome  $g_i$ , if we wish to include a hypothetical organism  $h$  that is 90% similar to genome  $g_i$  in its 30-mers, we can round down each entry of the column vector  $\text{pckm}_{30}(\cdot, g_i)$  to be no more than 0.90. Call this vector  $\text{pckm}_{30}(\cdot, h)$ . The entries below 90% do not need to be changed since we assume that the hypothetical organism has the same patterns of  $k$ -mer sharing to more distantly related “outgroup” taxa as to the reference organism.

We model the 50-mer similarity by setting  $\text{pckm}_{50}(\cdot, h) = p(\text{pckm}_{30}(\cdot, g_i))$  for  $p$  the previously defined polynomial. Adding these vectors to  $A^{(30)}$  and  $A^{(50)}$  effectively adds a hypothetical organism that has a common  $k$ -mer signature 90% similar to genome  $g_i$ . We then repeat this procedure for all training genomes  $g_i$  and for similarities ranging from 90%, 80%, ..., 10% and append these columns to  $A^{(30)}$  and  $A^{(50)}$ .

**2.3. Sample  $k$ -mer Signature.** Given a metagenomic sample, we form two vectors  $y^{(30)}$  and  $y^{(50)}$  consisting of the total counts in the sample of the 30-mers and 50-mers shared with the training organisms. In the Appendix section A.1, we show that these vectors are linearly related to the organism abundances via the common  $k$ -mer matrices  $A^{(30)}$  and  $A^{(50)}$ .

Note that in forming  $y^{(k)}$ , we count the  $k$ -mers in the entire sample, not of the individual reads. This allows for a very computationally efficient approach: as the training genomes typically have low error, their  $k$ -mers can be efficiently stored in de Bruijn graphs (formed using **Bcalm** [7]). We can then query the bloom count filter formed from the sample in a highly parallel fashion.

**2.4. Sparsity Promoting Optimization Procedure.** After forming  $y^{(k)}$ , we note that some of the entries  $y_i^{(k)}$  may be non-zero not due to the presence of organism  $i$  in the sample, but due to the fact that there exists an organism  $j$  that shares portions of its genome with organism  $i$ . Since  $A_{i,j}^{(k)}$  represents the “overlap” of these two organisms, we can deconvolute this linear mixture relationship by solving the equation  $A^{(k)}x = y^{(k)}$  for  $x$  the vector of organism abundances. However, after having augmented  $A^{(k)}$  with the hypothetical organisms, this system of equations is underdetermined (10 times more columns than rows). We can employ a sparsity promoting optimization procedure to infer the most parsimonious  $x$  consistent with the equations  $A^{(k)}x = y^{(k)}$  for  $k = 30, 50$ . This procedure, first introduced in [23] and proven correct in [12], is detailed in the Appendix section A.4.

**2.5. Inferring Taxonomy.** The abundances of the hypothetical organisms is then mapped back onto the taxonomy (for the output taxonomic profile) or the neighbor joining tree formed from the  $A^{(k)}$  (for the output strain variation figures) utilizing a least common ancestor approach detailed in the Appendix section A.5.

### 3. QUANTIFICATION OF STRAIN-LEVEL VARIATION

We demonstrate in two ways that the inclusion of the hypothetical organisms allows for the inference of strain-level variation. First, we spike novel organisms into a mock metagenomic community and show that **MetaPalette** can accurately predict their presence. Second, we utilize a real metagenomic soil sample to give evidence for a novel strain that **MetaPalette** predicts.

**3.1. HMP Mock Community.** We first formed the common  $k$ -mer matrices  $A^{(k)}$  using 31 strains of *Lysinibacillus sphaericus*. We then used **Grinder** [2] to simulate a dataset consisting of two novel strains (not included in the training database). These reads were then spiked into the HMP mock even community (a  $\sim$ 6.6M read metagenome consisting of 22 select organisms sampled using an Illumina GA-II sequencer; NCBI accession SRR172902). The output of **MetaPalette** is shown in Figure 5 demonstrating the ability of the method to correctly infer the presence of organisms absent from the training data.

Decreasing the number or changing the identity of the training organisms does not impede the method. In Figure 6, 50K simulated reads from the species *Providencia alcalifaciens* were again spiked into the HMP mock even community, and the inferred abundance is again placed optimally on the neighbor joining tree. Appendix section C contains a variety of such figures spanning all domains of life. These results provide evidence that MetaPalette can correctly infer the presence of organisms related to, but absent from, the training database.

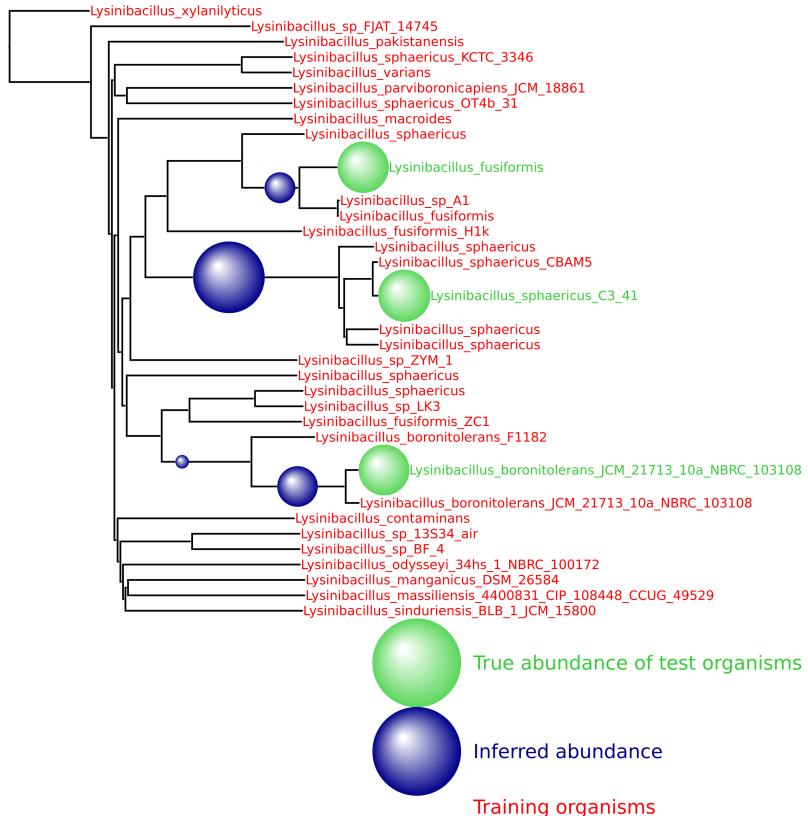


FIGURE 5. Result from training on 30 strains of *L. sphaericus* and testing on three novel strains. A total of 50K reads from the novel strains were spiked into the HMP mock even community. The training organisms are denoted with red names and the testing organisms have green names. The true abundance of the sample is pictured with green spheres and the inferred abundance is pictured with blue spheres.

**3.2. Metagenomic Soil Sample.** To assess MetaPalette on a real metagenomic sample, we utilized the Iowa prairie metagenomic sample from [18] (corresponding to MG-RAST project ID 6377). After running MetaPalette on a subset of this data (metagenome 4539594.3), the returned taxonomic profile predicted the presence of the genus *Bradyrhizobium*. Generating the tree plot on a subset of this genus resulted in, among others, a prediction of a novel organism in the clade defined by strains of *B. valentinum* (see Figure 7(a)). To verify this, we aligned the entire soil metagenome to the reference genome of the strain *B. valentinum LmjM3* using Bowtie2 with –very-sensitive-local settings [26] and extracted the aligned reads. Interestingly, 0.29% of the reads aligned, while the MetaPalette predicted abundance for this putative novel organism of interest was 0.33%. The depth of coverage of the extracted reads is pictured in Figure 7(b) and had a mean depth of 74.3X.

To assess the evolutionary relatedness of this predicted organism, we utilized the *B. valentinum LmjM3 nifH* gene sequence (NCBI accession KF806461) which was used in [10], along with other

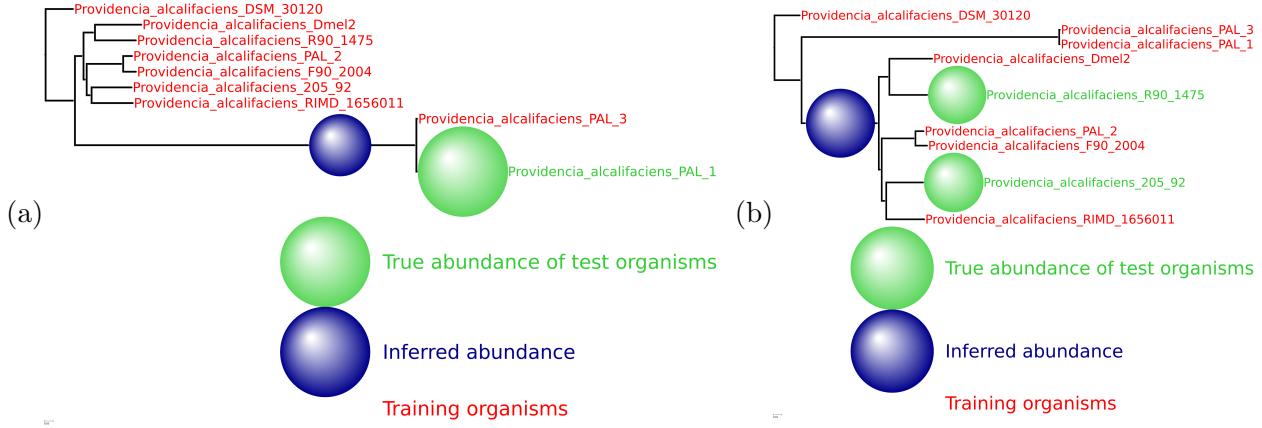


FIGURE 6. For each of the samples, a total of 50K reads from novel strains of *P. alcalifaciens* were spiked into the HMP mock even community. (a) Result from training on 8 strains and testing on one novel strain. (b) Result from training 7 strains and testing on two novel strains.

genes, to determine the taxonomy of *B. valentinum LmjM3*. Aligning the extracted reads to *nifH* resulted in a mean depth of coverage of 22X. We collapsed the aligned reads (via a majority vote) in regions of coverage at least 22X and called this the maximum likelihood sequence. We then performed a multiple sequence alignment of this sequence along with the *nifH* sequences of 20 other organisms closely related to *B. valentinum*. The topology of the bootstrap consensus neighbor-joining tree is pictured in Figure 7(c) and shows the maximum likelihood sequence is placed at the same location as was predicted by MetaPalette. While this is not enough evidence to unequivocally claim the existence of a novel strain in this sample, this gives support that MetaPalette correctly inferred the abundance and placement in Figure 7(a) of a potentially novel strain in the clade defined by strains of *B. valentinum*.

#### 4. COMPARISON TO OTHER METAGENOMIC PROFILING METHODS

To facilitate an objective comparison with other methods with minimal “author bias”, we utilized the same data and metrics used by other authors in a recent metagenomics methods evaluation paper [29]. This allowed comparison to the following algorithms: CLARK [37], Kraken [52], OneCodeX [35], LMAT [1], MG-RAST [34], MetaPhlAn [50], mOTU [47], Genometa [8], QIIME [5], EBI [19], MetaPhyler [30], MEGAN [20], taxator-tk [9], and GOTTCCHA [13].

**4.1. Training Data.** Each of the methods was trained using their default recommended databases. We trained our method using 6,914 whole genome sequences and assemblies obtained from various public repositories via RepoPhlAn (<https://bitbucket.org/nsegata/repophlan>). The training procedure for MetaPalette on these 6,914 organisms took a total of approximately 7 hours on a 48 core server.

**4.2. Testing Data.** The testing data consisted of 6 samples, and is fully explained in the Methods section of [29], but we briefly summarize it here. Three replicates were formed from two different distributions of over 900 different genomes spanning the tree of life (including Eukaryote genomes). Included in each test sample were shuffled/randomized genomes (not meant to be assigned to any known taxa) as well as sequences from the genome of *Leptospira interrogans* that were evolved using Rose [46] to simulate novelty. Error profiles were based on those of 6 real soil metagenomic samples sequenced using an Illumina HiSeq 2000. Each of the resulting test samples contains between 27 and 37 million read pairs.

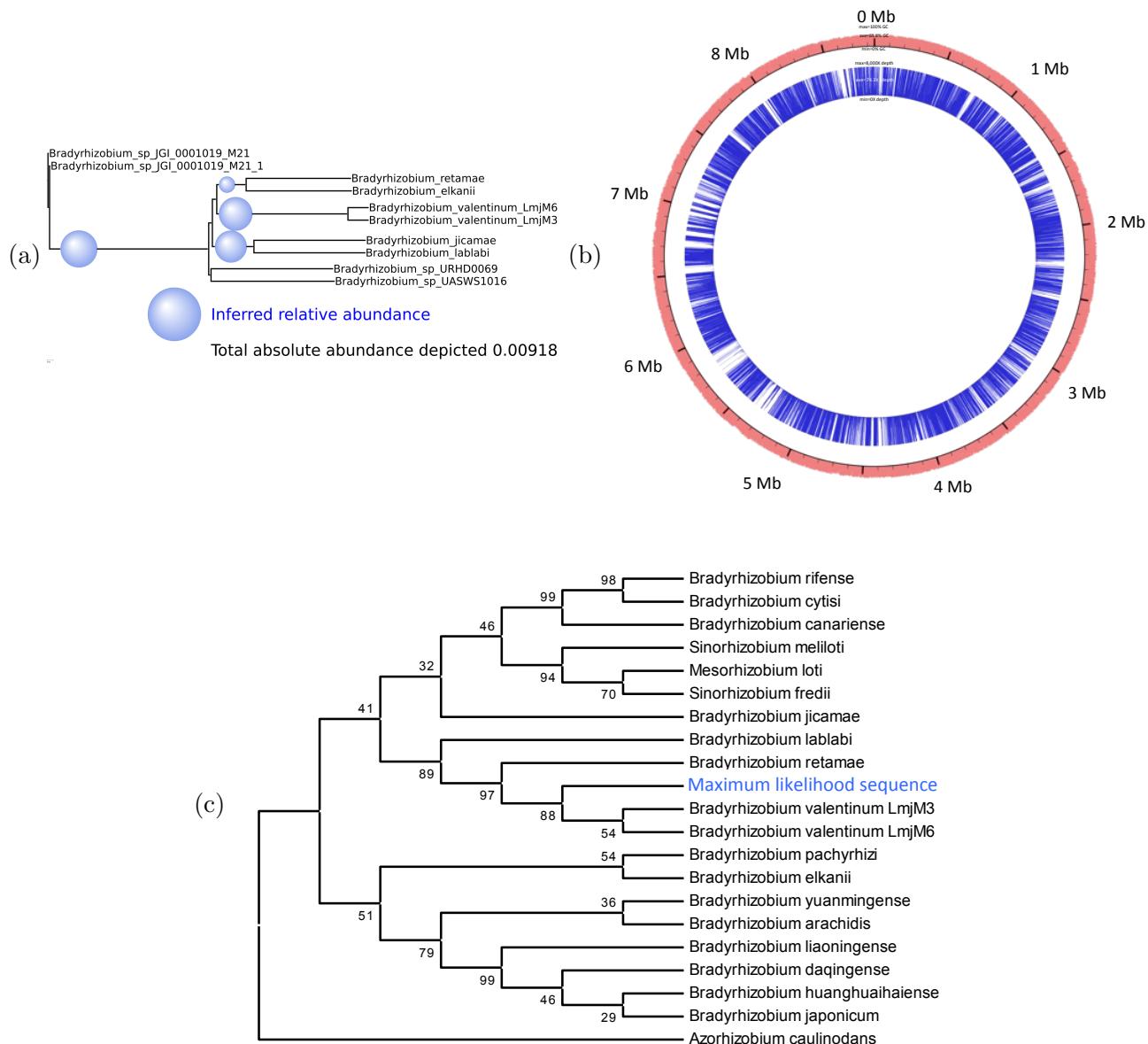


FIGURE 7. (a) Subtree of the MetaPalette output tree for the Iowa prairie metagenome using organisms from the genus *Bradyrhizobium*. (b) Depth of alignment for reads from the soil metagenome that aligned to *B. valentinum LmjM3*. The outer red ring shows the %GC for *B. valentinum LmjM3*, and the inner blue ring shows the alignment depth (truncated to 8,000X for ease of viewing). All contigs of the reference *B. valentinum LmjM3* were concatenated in this figure. (c) Bootstrap consensus tree topology based on *nifH* for 20 organisms along with the maximum likelihood sequence obtained from aligning the soil metagenome to the *nifH* gene sequence of *B. valentinum LmjM3*. Bootstrap values (500 replicates) are shown next to the branches. Full details regarding formation of the tree, along with a figure containing the branch lengths, are given in the Appendix section B.

**4.3. Error Metrics.** We utilized the same divergence error metric as [29], that is, for  $x_i$  representing the true frequency of taxa  $i$  in the sample, and  $x_i^*$  representing the predicted frequency for a given method of taxa  $i$ ,

$$\text{Divergence} = \sum_i \log_2 \left( \frac{x_i^*}{x_i} \right)$$

where the summation is over those indices such that  $x_i > 0$  and  $x_i^* > 0$ . Since this error metric does not take into consideration the number of spurious assignments (that is, taxa predicted by a method to be in a sample, but not actually present), we also use the number of false positives at a given taxonomic rank:

$$\text{FP} = |\{i : x_i^* > 0 \text{ and } x_i = 0\}|.$$

**4.4. Comparison Results.** Each method was run using the default parameters. For each method, we averaged the divergence error metric over all the test samples at the genus level (see Figure 8(a)). Furthermore, we selected a number of the more accurate methods and averaged the number of false positives over all the test samples at the phylum level (see Figure 8(b)). These two figures clearly show the competitive nature of MetaPalette as it has the lowest error in both metrics. However, when comparing to other methods, one should be careful of their intended use. For example, taxator-ik is intended to be used on an assembled metagenome (and here unassembled reads were used), and QIIME only uses the 16S rRNA sequences in a sample. Furthermore, most of these methods assign individual reads and then summarize this to obtain a taxonomic profile, while our method only profiles the entire sample and returns relative proportions of organisms.

Figure 9 shows the timing of each of the methods (on a log scale, obtained from [29]) further showing the competitive nature of MetaPalette.

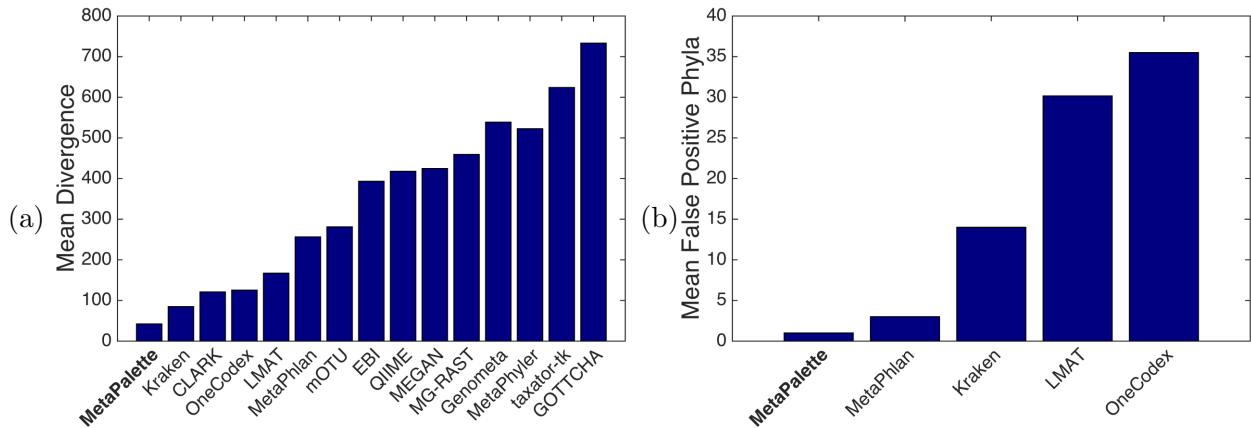


FIGURE 8. Plot of performance metrics for all methods averaged over all test samples. Smaller values indicate better performance. (a) Divergence error metric at the genus level. (b) Number of false positive phyla.

## 5. SOFTWARE AND PRE-TRAINED DATA

**5.1. Software.** The source code for MetaPalette, along with installation instructions and directions, is accessible at <https://github.com/dkoslicki/MetaPalette>. MetaPalette is written primarily in python and accepts input reads in uncompressed fasta or fastq format, as well as compressed fasta/fastq using bzip2 and gzip. For fastq input, optional parameters can be given to specify only counting  $k$ -mers above a certain quality score (Phred) thereby attenuating the negative impact of sequencing error in the correct inference of relative abundances. The output taxonomic profile is

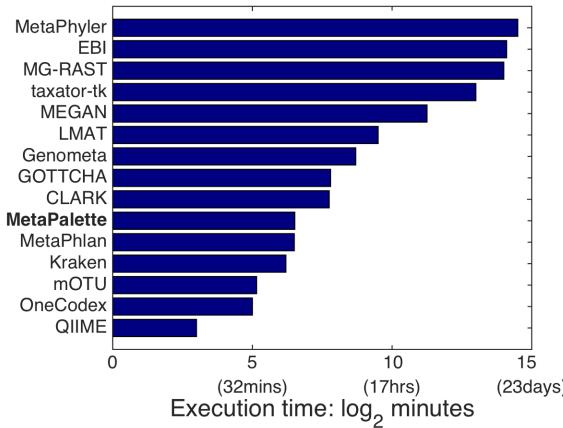


FIGURE 9. Mean execution time of each method averaged over all 6 test samples.

compliant with the Bioboxes profiling format version 0.9 found at <https://github.com/bioboxes/rfc/tree/master/data-format>. Python scripts are also included to aid in downloading data, forming custom databases, and creating the appropriate taxonomy files.

To facilitate cross-platform usability, a Docker [33] container has been created and is accessible at: <https://hub.docker.com/r/dkoslicki/metapalette> with an accompanying docker file at: <https://github.com/dkoslicki/MetaPalette/blob/master/Dockerfile>.

If users wish to use MetaPalette but lack computational resources, they may utilize the Galaxy [4, 14, 15] server located at: <http://math-galaxy.cgrb.oregonstate.edu/>.

A preliminary version of this software was submitted to the Critical Assessment of Metagenomic Interpretation (CAMI: <http://www.cami-challenge.org/>) under the name of CommonKmers. However, since significant changes have been made since that point, we strongly recommend using the current MetaPalette software instead.

**5.2. Pre-trained Data.** To decrease computational burden, pre-trained databases are accessible at [http://files.cgrb.oregonstate.edu/Koslicki\\_Lab/MetaPalette](http://files.cgrb.oregonstate.edu/Koslicki_Lab/MetaPalette). Databases and accompanying taxonomies have been included for archaea (666 organisms, 1.7 gigabytes uncompressed), bacteria (15,147 organisms, 60GB), eukaryota (1,307 organisms, 41GB), and viruses (4,798 organisms, 0.6GB). All organisms were obtained via RepoPhlAn.

The 6,914 organism database used for the comparison to other profiling methods is accessible at [http://files.cgrb.oregonstate.edu/Koslicki\\_Lab/MetaPallete/Comparison](http://files.cgrb.oregonstate.edu/Koslicki_Lab/MetaPallete/Comparison).

## 6. CONCLUSION

We have described a fast, flexible and accurate method for estimating the taxonomic composition of organisms which is based on reconstructing a  $k$ -mer-based profile of a sample. Each reference organism has a  $k$ -mer “palette” and we fit the sample as a mixture of different palettes, both of the reference organisms and organisms absent from the training data at varying degrees of relatedness to the training database. Our approach is in part inspired by chromosome painting method used to deduce fine-scale population structure in human genetics [17, 28] which is also based on mixture modeling of palettes. A particular advantage of MetaPalette over other metagenomic profiling methods is that MetaPalette provides an indication of how related the organisms in a given sample are to the closest matching organisms of the training database, whether they are within the same species or distantly related organisms from the same phyla.

Furthermore, the standard approach to summarizing composition information has been to place organisms at different taxonomic levels. We produce a standard taxonomic profile which we have

shown to be more accurate than that produced by other methods. This fixed rank approach is sensible at the genus level and above but omits fine-scale information. Hence, for branches of the tree of life that are well represented in the training database, we can also output a phylogenetic tree giving detailed information on how the sampled taxa relate to the organisms in the training database (Figures 5–7 and Appendix section C).

For many applications, it is of interest to understand which individual reads belong to which organisms [45]. A principled approach to this problem is to first estimate the overall composition of the sample, using MetaPalette or an equivalent, and then to assign individual reads conditional on the overall assignment. This represents a promising avenue for future methodological development.

## 7. ACKNOWLEDGEMENTS

The authors would like to thank the Isaac Newton Institute at Cambridge University for their hospitality during the program on metagnomics. This project was conceived while both authors were attending this program. Daniel Falush is supported by a Medical Research Fellow fellowship as part of the CLIMB Consortium for Medical Microbiology.

The authors would also like to thank Nam Nguyen and Daniel Alemany who both contributed to this project in its preliminary stages.

## APPENDIX A. TECHNICAL DETAILS OF THE COMMON $K$ -MER METHOD

**A.1. Mathematical Formulation.** We include here all rigorous mathematical definitions of the quantities discussed in the main text.

Given the alphabet  $\mathcal{A} = \{A, C, T, G\}$ , let  $\mathcal{A}^n$  denote the set of all words  $v$  of length  $|v| = n$  on  $\mathcal{A}$ , and let  $\mathcal{A}^* = \bigcup_{n \geq 0} \mathcal{A}^n$  be the set of all finite words on  $\mathcal{A}$ . Hence words containing non-*ACTG* characters are ignored. Let  $D = \{g_1, \dots, g_M\}$  be a database of genomic sequences  $g_j \in \mathcal{A}^*$  and let  $S = \{s_1, \dots, s_N\}$  be a set of sample sequences (the reads to be classified). For notational simplicity, assume that the read length is fixed: for all  $t$ ,  $|s_t| = r$ . Fix a  $k$ -mer size and endow  $\mathcal{A}^k = \{v_1, \dots, v_{4^k}\}$  with the lexicographic order. Let  $\text{occ}_v(w)$  represent the number of occurrences (with overlap) of the subword  $v$  in the word  $w$ . That is, for  $w, v \in \mathcal{A}^n$ , let

$$(A.1) \quad \text{occ}_v(w) = |\{j : w_j w_{j+1} \cdots w_{j+|v|-1} = v\}|.$$

For a fixed  $k$ -mer size, and two genomes  $g_i$  and  $g_j$  we calculate the number of  $k$ -mers in genome  $j$  common to both  $g_i$  and  $g_j$ . That is, the  $(i, j)^{\text{th}}$  entry of the common  $k$ -mer training matrix  $A^{(k)}$  is:

$$(A.2) \quad A_{i,j}^{(k)} = \sum_{w \in \text{SW}_{g_i}(k) \cap \text{SW}_{g_j}(k)} \frac{\text{occ}_w(g_j)}{|g_j| - k + 1}.$$

Refer to the entries of the common  $k$ -mer matrix as  $\text{pckm}_k(g_i, g_j) = A_{i,j}^{(k)}$ . Let  $s_t \subset g_j$  denote the relationship that read  $s_t$  was derived from genome  $g_j$ . We represent the taxonomic profile of the sample  $S$  by the probability vector  $x$ :

$$(A.3) \quad x_j = \frac{1}{N} \sum_{t=1}^N \mathbb{1}_{\{s_t \subset g_j\}}$$

where  $\mathbb{1}$  is the indicator function. Now let the measurement vector  $y$  be given by the probability vector

$$(A.4) \quad y_i^{(k)} = \sum_{w \in \text{SW}_{g_i}(k) \cap \text{SW}_S(k)} \text{occ}_w(S) / \sum_{w \in \text{SW}_S(k)} \text{occ}_w(S).$$

We assume that the reads  $s_t$  are uniformly randomly selected from the genomes  $g_j$ . Then for  $w \in \mathcal{A}^k$ , let  $\mathbb{P}(w|g_j)$  be the probability that  $k$ -mer  $w$  is found in genome  $g_j$ . Then we have that the

proportion of  $k$ -mers  $w$  in the sample is similar to the proportion of the appearance of  $w$  in the genomes  $g_j$  when weighted by the relative abundance of the genomes  $g_i$  in the sample:

$$(A.5) \quad \frac{\sum_{t=1}^N \text{occ}_w(s_t)}{\sum_{w \in \text{SW}_S(k)} \text{occ}_w(S)} \approx \frac{1}{N} \sum_{j=1}^M \mathbb{P}(w|g_j) \sum_{t=1}^N \mathbb{1}_{\{s_t \subset g_j\}}$$

$$(A.6) \quad = \frac{1}{N} \sum_{j=1}^M \frac{\text{occ}_{g_j}(w)}{|g_j| - n + 1} \sum_{t=1}^N \mathbb{1}_{\{s_t \subset g_j\}}$$

We then calculate

$$(A.7) \quad y_i^{(k)} = \frac{1}{\sum_{w \in \text{SW}_S(k)} \text{occ}_w(S)} \sum_{w \in \text{SW}_{g_i}(k) \cap \text{SW}_S(k)} \text{occ}_w(S)$$

$$(A.8) \quad = \frac{1}{\sum_{w \in \text{SW}_S(k)} \text{occ}_w(S)} \sum_{w \in \text{SW}_{g_i}(k) \cap \text{SW}_S(k)} \sum_{t=1}^N \text{occ}_w(s_t)$$

$$(A.9) \quad \approx \frac{1}{N} \sum_{w \in \text{SW}_{g_i}(k) \cap \text{SW}_S(k)} \sum_{j=1}^M \frac{\text{occ}_{g_j}(w)}{|g_j| - k + 1} \sum_{t=1}^N \mathbb{1}_{\{s_t \subset g_j\}}$$

$$(A.10) \quad = \frac{1}{N} \sum_{j=1}^M \sum_{w \in \text{SW}_{g_i}(k) \cap \text{SW}_{g_j}(k) \cap \text{SW}_S(k)} \frac{\text{occ}_{g_j}(w)}{|g_j| - k + 1} \sum_{t=1}^N \mathbb{1}_{\{s_t \subset g_j\}}$$

$$(A.11) \quad \approx \frac{1}{N} \sum_{j=1}^M \sum_{w \in \text{SW}_{g_i}(k) \cap \text{SW}_{g_j}(k)} \frac{\text{occ}_{g_j}(w)}{|g_j| - k + 1} \sum_{t=1}^N \mathbb{1}_{\{s_t \subset g_j\}}$$

$$(A.12) \quad = \sum_{j=1}^M A_{i,j}^{(k)} x_j = (A^{(k)} x)_i$$

Line (A.10) is justified since if  $w \notin \text{SW}_{g_j}(k)$  then  $\text{occ}_{g_j}(w) = 0$ . For computational reasons, we make the assumption in line (A.11) that  $\text{SW}_{g_i}(k) \cap \text{SW}_{g_j}(k) \cap \text{SW}_S(k) = \text{SW}_{g_i}(k) \cap \text{SW}_{g_j}(k)$ . However, this assumption can be mitigated by adding hypothetical organisms (see section A.3). Our assumptions imply that

$$(A.13) \quad A^{(k)} x \approx y^{(k)}.$$

We will try to recover the vector  $x$  satisfying  $x_j \geq 0$  for all  $j = 1, \dots, M$  from equation (A.13).

**A.2. Further Improvements.** A few further improvements are possible, but not pursued here. Namely, we could use just the  $k$ -mers that are actually in the sample to form the training matrix. I.e use  $\text{SW}_{g_i}(k) \cap \text{SW}_{g_j}(k) \cap \text{SW}_S(k)$  in the formation of  $A^{(k)}$ :

$$A_{i,j}^{(k)} = \sum_{w \in \text{SW}_{g_i}(k) \cap \text{SW}_{g_j}(k) \cap \text{SW}_S(k)} \frac{\text{occ}_w(g_j)}{|g_j| - k + 1}.$$

The disadvantage of this is that the (slow) training step would need to be re-run for each sample.

For a second improvement, we could make the approximation in (A.5) more delicate by incorporating the coverage:

$$(A.14) \quad \frac{\sum_{t=1}^N \text{occ}_w(s_t)}{\sum_{w \in \text{SW}_S(k)} \text{occ}_w(S)} \approx \frac{1}{N} \sum_{j=1}^M \mathbb{P}(w|g_j) \frac{\sum_{w \in \text{SW}_{g_j}(k) \cap \text{SW}_S(k)} \text{occ}_w(g_j)}{\sum_{w \in \text{SW}_{g_j}(k)} \text{occ}_w(g_j)} \sum_{t=1}^N \mathbb{1}_{\{s_t \subset g_j\}}$$

$$(A.15) \quad = \frac{1}{N} \sum_{j=1}^M \frac{\text{occ}_{g_j}(w)}{|g_j| - n + 1} \frac{\sum_{w \in \text{SW}_{g_j}(k) \cap \text{SW}_S(k)} \text{occ}_w(g_j)}{|g_j| - k + 1} \sum_{t=1}^N \mathbb{1}_{\{s_t \subset g_j\}}.$$

So  $A^{(k)}$  would have the form:

$$(A.16) \quad A_{i,j}^{(k)} = \frac{\sum_{w \in \text{SW}_{g_j}(k) \cap \text{SW}_S(k)} \text{occ}_w(g_j)}{|g_j| - k + 1} \sum_{w \in \text{SW}_{g_i}(k) \cap \text{SW}_{g_j}(k) \cap \text{SW}_S(k)} \frac{\text{occ}_w(g_j)}{|g_j| - k + 1}.$$

This effectively multiplies the columns  $j$  of  $A^{(k)}$  by the % coverage of genome  $j$ . Lastly, in (A.16), we could put a weighting factor that represents how unique a  $k$ -mer is to the genome in question. This would down-weight  $k$ -mers shared among many diverse genomes and up-weight those unique to certain strains/species/genera/etc.

**A.3. Hypothetical Organisms.** To simulate an organism that is, say, 90% related to a database genome  $g_i$ , we augment the common  $k$ -mer matrix  $A^{(k)}$  with a column derived by rounding down the entries of the column vector  $\{A_{i,j}^{(k)}\}_{i=1,\dots,M}$  that are above 90%. Two  $k$ -mer sizes are needed to form the hypothetical organism common  $k$ -mer matrices. For the first  $k$ -mer size,  $k_1$ , we define  $A^{(k_1),h}$  for a fixed number of hypothetical bins  $h \in \{0.9, 0.8, \dots, 0.1\}$  where

$$(A.17) \quad A_{i,j}^{(k_1),h} = \max \left( h, \sum_{w \in \text{SW}_{g_i}(k_1) \cap \text{SW}_{g_j}(k_1)} \frac{\text{occ}_w(g_j)}{|g_j| - n + 1} \right).$$

For the second  $k$ -mer size,  $k_2$ , using the polynomial  $p(x) = -.5141x^3 + 1.0932x^2 + 0.3824x$ , we define

$$(A.18) \quad A_{i,j}^{(k_2),h} = \max \left( h, p \left( \sum_{w \in \text{SW}_{g_i}(k_2) \cap \text{SW}_{g_j}(k_2)} \frac{\text{occ}_w(g_j)}{|g_j| - n + 1} \right) \right).$$

Instead of thresholding, as we did here, one can imagine other scalings obtained from studying the relationship between a given taxonomy and the common  $k$ -mer matrix  $A^{(k)}$ . In particular, to deal with differing rates of evolution in the tree of life, a fruitful area of future investigation would be to modify the polynomial  $p$  depending on the taxonomy of the organisms under consideration.

**A.4. Optimization Procedure.** We choose two  $k$ -mer sizes to be  $k = 30$  and  $k = 50$ , as this seems to give a good trade-off between reconstruction fidelity and computational performance. We then collect the common  $k$ -mer matrix and hypothetical matrices block-wise into the  $2|D| \times 10|D|$  size matrix

$$(A.19) \quad A = \begin{bmatrix} A^{(30)}, & A^{(30),0.9}, & \dots, & A^{(30),0.1} \\ A^{(50)}, & A^{(50),0.9}, & \dots, & A^{(50),0.1} \end{bmatrix}.$$

Collect also the  $k$ -mer sample vectors  $y^{(k)}$ :

$$(A.20) \quad y = \begin{bmatrix} y^{(30)} \\ y^{(50)} \end{bmatrix}.$$

The problem at hand is then to reconstruct the phylogenetic profile  $x$  by solving the linear system

$$(A.21) \quad Ax = y.$$

Equation (A.21) is solved by using a sparsity-promoting optimization procedure motivated by techniques used in the compressive sensing literature. Sparsity is emphasized due to the inclusion of the hypothetical organisms, as well as the reasonable assumption that relatively few organisms from the database  $D$  are actually present in the given sample. We use a variant of nonnegative basis pursuit denoising which reduces to a nonnegative least squares problem [6, 12]. We aim to solve

$$(\ell_1\text{-min}) \quad x^* = \arg \min_z \|z\|_1 \quad \text{subject to } Az = y, z \geq 0.$$

This optimization procedure has the advantage of being transformed into a nonnegative least squares problem. Indeed as  $\lambda \rightarrow \infty$ , we can regularize  $(\ell_1\text{-min})$  as

$$(\text{NNREG}) \quad x^* = \arg \min_z \|z\|_1^2 + \lambda^2 \|Az - y\|_2^2 \quad \text{subject to } z \geq 0.$$

This reduces to a nonnegative least squares problem by defining

$$\tilde{A} := \left[ \frac{1 \cdots 1}{\lambda A} \right], \quad \tilde{y} := \left[ \frac{0}{\lambda y} \right].$$

So (NNREG) is equivalent to the nonnegative least squares problem

$$x^* = \arg \min_z \|\tilde{A}z - \tilde{y}\|_2^2 \quad \text{subject to } z \geq 0.$$

This can be solved efficiently by using the Lawson–Hanson algorithm [27]. We use the value  $\lambda = 200$  throughout as this value gives a good trade-off between sparsity and accuracy of fit of the  $k$ -mer counts.

**A.5. Inferring Taxonomy.** Since the reconstructed vector  $x$  may have non-zero entries corresponding to a hypothetical bin, we need to develop a method to map from a hypothetical bin to a specific taxonomic rank. A naïve approach would be to assign a fixed taxonomic rank to each hypothetical bin (call this the *fixed rank* method). For example, all non-zero entries of  $x$  corresponding to  $A^{(k)}$  would be assigned to the strain level, all non-zero entries of  $x$  corresponding to  $A^{(k),0.9}$  would be assigned to the species level, etc.

We take a more biologically informed approach: we take the least common ancestor (LCA) taxa between a hypothetical organism and a nearby organism in the database  $D$ : if  $x_i > 0$  corresponds to the hypothetical bin  $h$ , find an organism  $g_j$  such that  $|A_{i,j}^{(k)} - h| < \delta$  for some threshold  $\delta$ . In the output taxonomic profile, we assign  $x_i$  to the lowest taxonomic rank common to organisms  $g_i$  and  $g_j$ . For the output strain variation figures, we assign the abundance  $x_i$  relative to the least common ancestor of  $g_i$  and  $g_j$  (above the LCA if  $h < A_{i,j}^{(k)}$  and below the LCA if  $h > A_{i,j}^{(k)}$ ).

For the output taxonomic profile, a hybrid of the fixed rank and LCA approaches can increase sensitivity or specificity. We thus include three options: the default option is the LCA approach, while the sensitive and specific options are varying hybrids of the two methods.

## APPENDIX B. SEQUENCE ANALYSIS DETAILS

In Figure 7(c), the evolutionary history was inferred using the Neighbor-Joining method [42]. The bootstrap consensus tree inferred from 500 replicates is taken to represent the evolutionary history of the taxa analyzed [11]. Branches corresponding to partitions reproduced in less than 50% bootstrap replicates are collapsed. The percentage of replicate trees in which the associated taxa clustered together in the bootstrap test (500 replicates) are shown next to the branches. The analysis involved 21 nucleotide sequences. All positions with less than 95% site coverage were eliminated. That is, fewer than 5% alignment gaps, missing data, and ambiguous bases were allowed at any position. There were a total of 652 positions in the final dataset. Evolutionary analyses were conducted in MEGA6 [49].

Figure 10 depicts a tree using the same method as just described, but with evolutionary distances computed using the Maximum Composite Likelihood method [48]. Units are the number of base substitutions per site.

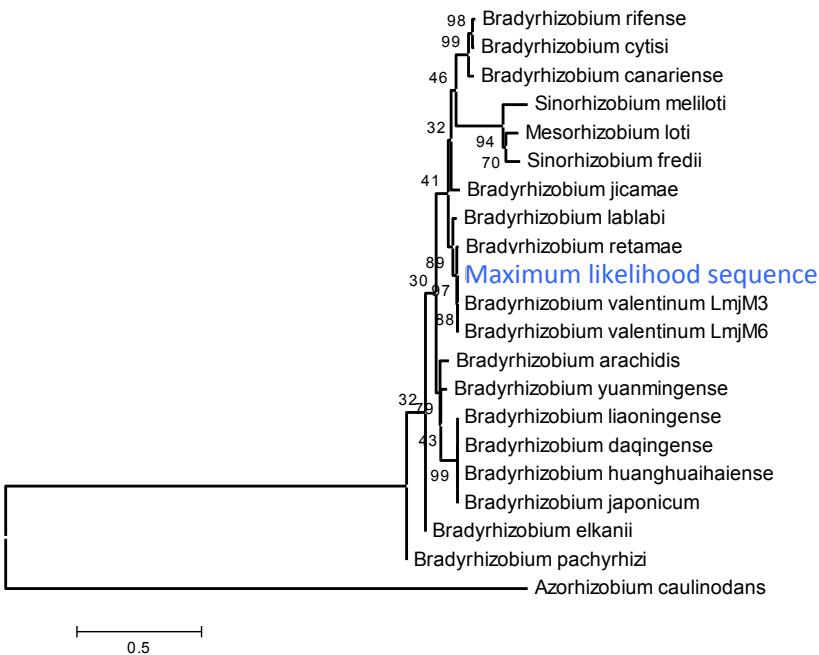


FIGURE 10. Neighborhood-joining tree based on *nifH* for 20 organisms along with the maximum likelihood sequence obtained from aligning the soil data to the gene sequence for *nifH* of *B. valentinum* *LmjM3*. Bootstrap values are shown next to the branches and the bar indicates 0.5 nucleotide substitutions per site.

#### APPENDIX C. ADDITIONAL FIGURES

We provide here a number of additional output figures from MetaPalette to demonstrate that the ability to correctly infer the presence of organisms related to, but absent from, the training database is not dependent on the particular kingdom/phyla/etc. used. Unless otherwise noted, a total of 50K simulated reads from the novel organisms were spiked into the HMP mock even community. Figures are included for Bacteria, Archaea, Eukaryota, and viruses.

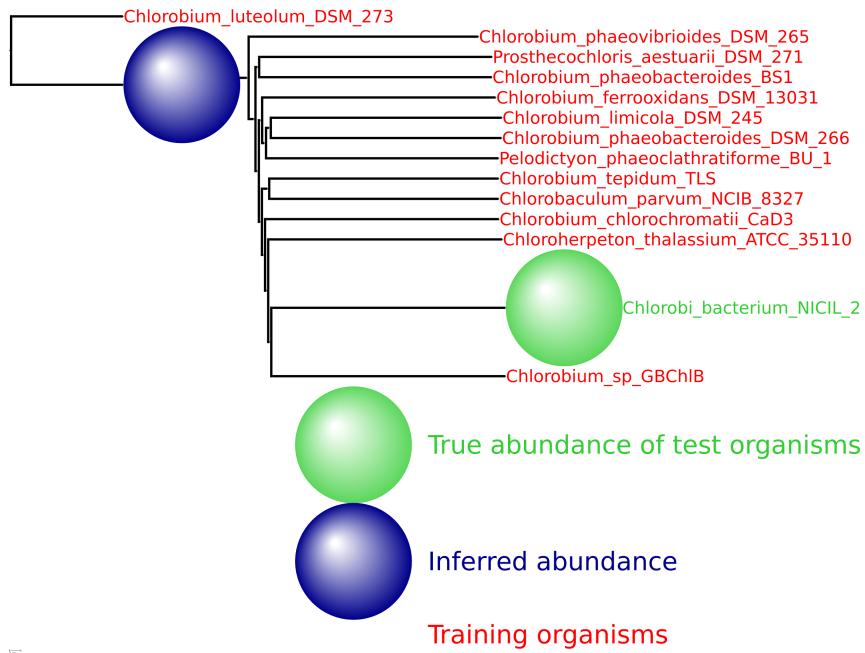


FIGURE 11. HMP spiking results for the bacterial phylum *Chlorobi* with 13 training organisms and one novel testing organism.

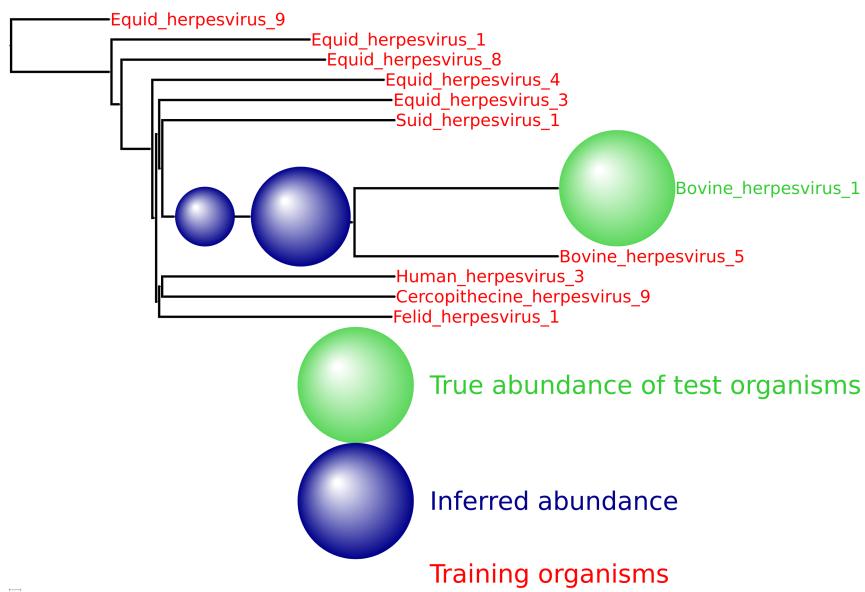


FIGURE 12. HMP spiking results for the viral genus *Varicellovirus* with 10 training organism and one novel testing organism. Only 5K simulated reads were spiked into the HMP sample.

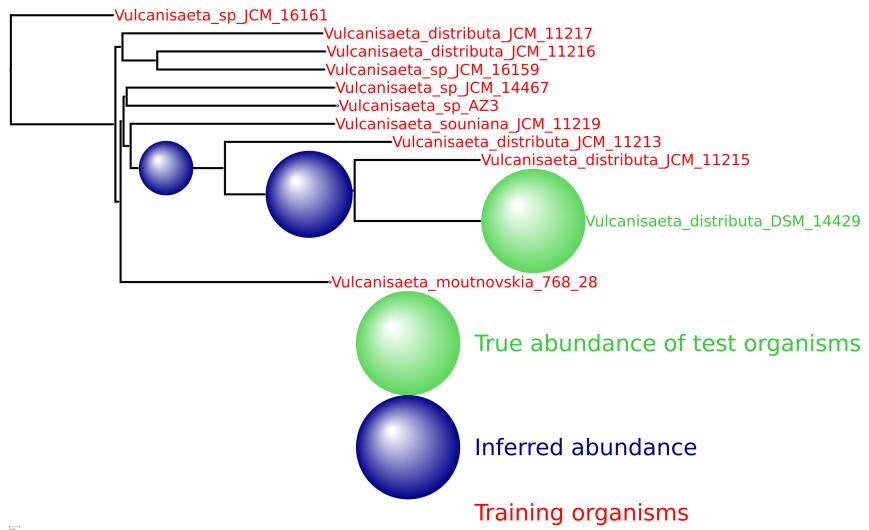


FIGURE 13. HMP spiking results for the archaeal genus *Vulcanisaeta* with 10 training organism and one novel testing organism. Only 5K simulated reads were spiked into the HMP sample.

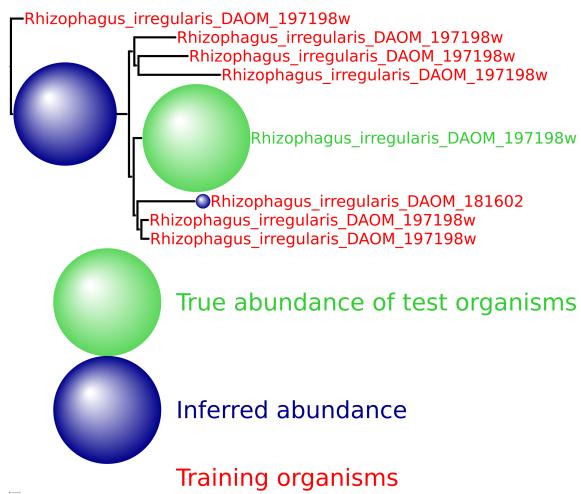


FIGURE 14. HMP spiking results for the eukaryota genus *Rhizophagus* with 7 training organism and one novel testing organism.

## REFERENCES

- [1] S. K. Ames, D. A. Hysom, S. N. Gardner, G. S. Lloyd, M. B. Gokhale, and J. E. Allen. Scalable metagenomic taxonomy classification using a reference genome database. *Bioinformatics*, 29(18):2253–2260, 2013.
- [2] F. E. Angly, D. Willner, F. Rohwer, P. Hugenholtz, and G. W. Tyson. Grinder: a versatile amplicon and shotgun sequence simulator. *Nucleic acids research*, 61(0):1–8, Mar. 2012.
- [3] B. J. Baker, L. R. Comolli, G. J. Dick, L. J. Hauser, D. Hyatt, B. D. Dill, M. L. Land, N. C. VerBerkmoes, R. L. Hettich, and J. F. Banfield. Enigmatic, ultrasmall, uncultivated archaea. *Proceedings of the National Academy of Sciences*, 107(19):8806–8811, 2010.
- [4] D. Blankenberg, G. V. Kuster, N. Coraor, G. Ananda, R. Lazarus, M. Mangan, A. Nekrutenko, and J. Taylor. Galaxy: a web-based genome analysis tool for experimentalists. *Current protocols in molecular biology*, pages 19–10, 2010.
- [5] J. G. Caporaso, J. Kuczynski, J. Stombaugh, K. Bittinger, F. D. Bushman, E. K. Costello, N. Fierer, A. G. Peña, J. K. Goodrich, J. I. Gordon, et al. Qiime allows analysis of high-throughput community sequencing data. *Nature methods*, 7(5):335–336, 2010.
- [6] S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic Decomposition by Basis Pursuit. *SIAM Journal on Scientific Computing*, 20(1):33–61, Jan. 1998.
- [7] R. Chikhi, A. Limasset, S. Jackman, J. T. Simpson, and P. Medvedev. On the representation of de bruijn graphs. In *Research in Computational Molecular Biology*, pages 35–55. Springer, 2014.
- [8] C. F. Davenport, J. Neugebauer, N. Beckmann, B. Friedrich, B. Kameri, S. Kokott, M. Paetow, B. Siekmann, M. Wieding-Drewes, M. Wienhfer, S. Wolf, B. Tmmler, V. Ahlers, and F. Sprengel. Genometa - a fast and accurate classifier for short metagenomic shotgun reads. *PLoS ONE*, 7(8):e41224, 08 2012.
- [9] J. Dröge, I. Gregor, and A. McHardy. Taxator-tk: precise taxonomic assignment of metagenomes by fast approximation of evolutionary neighborhoods. *Bioinformatics*, 31(6):817–824, 2015.
- [10] D. Durán, L. Rey, A. Busquets, J. Imperial, and T. Ruiz-Argüeso. *Bradyrhizobium valentinum* sp. nov., isolated from effective nodules of *lupinus mariae-josephae*, a lupine endemic of basic-lime soils in eastern spain. *Systematic and applied microbiology*, 37(5):336–341, 2014.
- [11] J. Felsenstein. Confidence limits on phylogenies: an approach using the bootstrap. *Evolution*, pages 783–791, 1985.
- [12] S. Foucart and D. Koslicki. Sparse recovery by means of nonnegative least squares. *IEEE Signal Processing Letters*, 21(4):498–502, 2014.
- [13] T. A. K. Freitas, P.-E. Li, M. B. Scholz, and P. S. Chain. Accurate read-based metagenome characterization using a hierarchical suite of unique signatures. *Nucleic acids research*, page gkv180, 2015.
- [14] B. Giardine, C. Riemer, R. C. Hardison, R. Burhans, L. Elnitski, P. Shah, Y. Zhang, D. Blankenberg, I. Albert, J. Taylor, et al. Galaxy: a platform for interactive large-scale genome analysis. *Genome research*, 15(10):1451–1455, 2005.
- [15] J. Goecks, A. Nekrutenko, J. Taylor, et al. Galaxy: a comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences. *Genome Biol*, 11(8):R86, 2010.
- [16] D. H. Haft and A. Tovchigrechko. High-speed microbial community profiling. *Nature methods*, 9(8):793–794, 2012.
- [17] G. Hellenthal, G. B. Busby, G. Band, J. F. Wilson, C. Capelli, D. Falush, and S. Myers. A genetic atlas of human admixture history. *Science*, 343(6172):747–751, 2014.
- [18] A. C. Howe, J. K. Jansson, S. A. Malfatti, S. G. Tringe, J. M. Tiedje, and C. T. Brown. Tackling soil diversity with the assembly of large, complex metagenomes. *Proceedings of the National Academy of Sciences*, 111(13):4904–4909, 2014.
- [19] S. Hunter, M. Corbett, H. Denise, M. Fraser, A. Gonzalez-Beltran, C. Hunter, P. Jones, R. Leinonen, C. McAnulla, E. Maguire, et al. Ebi metagenomicsa new resource for the analysis and archiving of metagenomic data. *Nucleic acids research*, 42(D1):D600–D606, 2014.
- [20] D. H. Huson, S. Mitra, H.-J. Ruscheweyh, N. Weber, and S. C. Schuster. Integrative analysis of environmental sequences using megan4. *Genome research*, 21(9):1552–1560, 2011.
- [21] P. Jaccard. *Etude comparative de la distribution florale dans une portion des Alpes et du Jura*. Impr. Corbaz, 1901.
- [22] J. B. Kaper, J. P. Nataro, and H. L. Mobley. Pathogenic escherichia coli. *Nature Reviews Microbiology*, 2(2):123–140, 2004.
- [23] D. Koslicki, S. Foucart, and G. Rosen. Quikr: a method for rapid reconstruction of bacterial communities via compressive sensing. *Bioinformatics*, page btt336, 2013.
- [24] D. Koslicki, S. Foucart, and G. Rosen. WGSQuikr: fast whole-genome shotgun metagenomic classification. *PloS one*, 9(3):91784, 2014.

- [25] Y. Lan, Q. Wang, J. R. Cole, and G. L. Rosen. Using the rdp classifier to predict taxonomic novelty and reduce the search space for finding novel organisms. *PLoS one*, 7(3):e32491, 2012.
- [26] B. Langmead and S. L. Salzberg. Fast gapped-read alignment with bowtie 2. *Nature methods*, 9(4):357–359, 2012.
- [27] C. L. Lawson and R. J. Hanson. *Solving least squares problems*, volume 161. SIAM, 1974.
- [28] S. Leslie, B. Winney, G. Hellenthal, D. Davison, A. Boumertit, T. Day, K. Hutnik, E. C. Rorvik, B. Cunliffe, D. J. Lawson, et al. The fine-scale genetic structure of the british population. *Nature*, 519(7543):309–314, 2015.
- [29] S. Lindgreen, K. L. Adair, and P. Gardner. An evaluation of the accuracy and speed of metagenome analysis tools. *bioRxiv*, page 017830, 2015.
- [30] B. Liu, T. Gibbons, M. Ghodsi, and M. Pop. Metaphyler: Taxonomic profiling for metagenomic sequences. In *Bioinformatics and Biomedicine (BIBM), 2010 IEEE International Conference on*, pages 95–100. IEEE, 2010.
- [31] G. Marçais and C. Kingsford. A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics*, 27(6):764–770, 2011.
- [32] Y. Marcy, C. Ouverney, E. M. Bik, T. Lösekann, N. Ivanova, H. G. Martin, E. Szeto, D. Platt, P. Hugenholtz, D. A. Relman, et al. Dissecting biological dark matter with single-cell genetic analysis of rare and uncultivated tm7 microbes from the human mouth. *Proceedings of the National Academy of Sciences*, 104(29):11889–11894, 2007.
- [33] D. Merkel. Docker: lightweight linux containers for consistent development and deployment. *Linux Journal*, 2014(239):2, 2014.
- [34] F. Meyer, D. Paarmann, M. D’Souza, R. Olson, E. M. Glass, M. Kubal, T. Paczian, A. Rodriguez, R. Stevens, A. Wilke, et al. The metagenomics rast server—a public resource for the automatic phylogenetic and functional analysis of metagenomes. *BMC bioinformatics*, 9(1):386, 2008.
- [35] S. S. Minot, N. Krumm, and N. B. Greenfield. One codex: A sensitive and accurate data platform for genomic microbial identification. *bioRxiv*, page 027607, 2015.
- [36] N.-p. Nguyen, S. Mirarab, B. Liu, M. Pop, and T. Warnow. Tipp: taxonomic identification and phylogenetic profiling. *Bioinformatics*, 30(24):3548–3555, 2014.
- [37] R. Ounit, S. Wanamaker, T. J. Close, and S. Lonardi. Clark: fast and accurate classification of metagenomic and genomic sequences using discriminative k-mers. *BMC genomics*, 16(1):236, 2015.
- [38] H. Philippe, H. Brinkmann, D. V. Lavrov, D. T. J. Littlewood, M. Manuel, G. Wörheide, and D. Baurain. Resolving difficult phylogenetic questions: why more sequences are not enough. *PLoS Biol*, 9(3):e1000602, 2011.
- [39] P. Puigbò, Y. I. Wolf, and E. V. Koonin. Search for a’tree of life’ in the thicket of the phylogenetic forest. *Journal of Biology*, 8(6):1, 2009.
- [40] C. Rinke, P. Schwientek, A. Sczyrba, N. N. Ivanova, I. J. Anderson, J.-F. Cheng, A. Darling, S. Malfatti, B. K. Swan, E. A. Gies, et al. Insights into the phylogeny and coding potential of microbial dark matter. *Nature*, 499(7459):431–437, 2013.
- [41] B. Roura, D. Baurain, and H. Philippe. Impact of missing data on phylogenies inferred from empirical phylogenomic data sets. *Molecular biology and evolution*, 30(1):197–214, 2013.
- [42] N. Saitou and M. Nei. The neighbor-joining method: a new method for reconstructing phylogenetic trees. *Molecular biology and evolution*, 4(4):406–425, 1987.
- [43] L. Salichos and A. Rokas. Inferring ancient divergences requires genes with strong phylogenetic signals. *Nature*, 497(7449):327–331, 2013.
- [44] L. Schaeffer, H. Pimentel, N. Bray, P. Melsted, and L. Pachter. Pseudoalignment for metagenomic read assignment. *arXiv preprint arXiv:1510.07371*, 2015.
- [45] T. J. Sharpton. An introduction to the analysis of shotgun metagenomic data. *Frontiers in plant science*, 5, 2014.
- [46] J. Stoye, D. Evers, and F. Meyer. Rose: generating sequence families. *Bioinformatics*, 14(2):157–163, 1998.
- [47] S. Sunagawa, D. R. Mende, G. Zeller, F. Izquierdo-Carrasco, S. A. Berger, J. R. Kultima, L. P. Coelho, M. Arumugam, J. Tap, H. B. Nielsen, et al. Metagenomic species profiling using universal phylogenetic marker genes. *Nature methods*, 10(12):1196–1199, 2013.
- [48] K. Tamura, M. Nei, and S. Kumar. Prospects for inferring very large phylogenies by using the neighbor-joining method. *Proceedings of the National Academy of Sciences of the United States of America*, 101(30):11030–11035, 2004.
- [49] K. Tamura, G. Stecher, D. Peterson, A. Filipski, and S. Kumar. Mega6: molecular evolutionary genetics analysis version 6.0. *Molecular biology and evolution*, page mst197, 2013.
- [50] D. T. Truong, E. A. Franzosa, T. L. Tickle, M. Scholz, G. Weingart, E. Pasolli, A. Tett, C. Huttenhower, and N. Segata. Metaphlan2 for enhanced metagenomic taxonomic profiling. *Nature methods*, 12(10):902–903, 2015.
- [51] B. Wang. Limitations of compositional approach to identifying horizontally transferred genes. *Journal of molecular evolution*, 53(3):244–250, 2001.
- [52] D. E. Wood and S. L. Salzberg. Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome Biol*, 15(3):R46, 2014.