

MicroPheno: Predicting environments and host phenotypes from 16S rRNA gene sequencing using a k-mer based representation of shallow sub-samples

Ehsaneddin Asgari^{1,3}, Kiavash Garakani², Alice C. McHardy³, and Mohammad R.K. Mofrad^{1,2,*}

¹Department of Bioengineering, University of California, Berkeley, CA 94720, USA

²Molecular Biophysics and Integrated Bioimaging, Lawrence Berkeley National Lab, Berkeley, CA 94720, USA

³Computational Biology of Infection Research, Helmholtz Centre for Infection Research, Brunswick 38124, Germany

*mofrad@berkeley.edu

ABSTRACT

Motivation: Microbial communities play important roles in the function and maintenance of various biosystems, ranging from the human body to the environment. A major challenge in microbiome research is the classification of microbial communities of different environments or host phenotypes. The most common and cost-effective approach for such studies to date is 16S rRNA gene sequencing. Recent falls in sequencing costs have increased the demand for simple, efficient, and accurate methods for rapid detection or diagnosis with proved applications in medicine, agriculture, and forensic science.

Results: We describe a reference- and alignment-free approach for predicting the environment or host phenotype from microbial community samples based on k-mer distributions in 16S rRNA data. In addition, we propose a bootstrapping framework to investigate the sufficiency of a shallow sub-sample for prediction. We study the use of deep learning methods as well as classic machine learning approaches for distinguishing among human body-sites, diagnosis of Crohn's disease, and predicting the environments (18 ecological and 5 organismal environments) from representative 16S gene sequences. Furthermore, we explore the use of unsupervised dimensionality reduction methods as well as supervised deep representation learning for visualizing microbial data of different environments and host phenotypes. We demonstrated that k-mer representations outperform Operational Taxonomic Unit (OTU) features in distinguishing among 5 major body-sites, as well as predicting Crohn's disease using 16S rRNA sequencing samples. We also showed that a shallow sub-sample of 16S rRNA samples alone can be sufficient to produce a proper k-mer representation of data. Aside from being more accurate, using k-mer features in shallow sub-samples provided the following benefits: (i) skipping computationally costly sequence alignments required in OTU-picking, and (ii) proof of concept for the sufficiency of a shallow and short-length 16S rRNA sequencing for environment/host phenotype prediction. In addition, k-mer features were able to accurately predict representative sequences of 18 ecological and 5 organismal environments with relatively high macro-F1 scores. Deep Neural Network outperformed Random Forest and Support Vector Machine in classification of large datasets.

Availability: The link to the MicroPheno code and the datasets will be available at <https://llp.berkeley.edu/micropheno>.

Frequent abbreviations used are \mathbf{D}_{KL} : Kullback Leibler divergence, $\mathbf{DNN-l}$: Deep Neural Network with l layers, \mathbf{RF} : Random Forests, and \mathbf{SVM} : Support Vector Machine (here linear SVM).

1 Introduction

Microbial communities have important functions relevant to supporting, regulating, and in some cases causing unwanted conditions (e.g. diseases or pollution) in their hosts/environments, ranging from organismal environments, like the human body, to ecological environments, like soil and water. These communities typically consist of a variety of microorganisms, including eukaryotes, archaea, bacteria and viruses. The human microbiome is the set of all microorganisms that live in close association with the human body. It is now widely believed that changes in the composition of our microbiomes correlate with numerous disease states, raising the possibility that manipulation of these communities may be used to treat diseases. Normally, the microbiota (particularly the intestinal microbiota) of the microbiome play several important roles in humans, which include: (i) prevention of pathogen growth, (ii) education and regulation of the host immune system, and (iii) providing energy substrates to the host [1]. Consequently, dysbiosis of the human microbiome can promote several diseases, including asthma [2, 3], irritable bowel syndrome [4, 5], *Clostridium difficile* infection [6], chronic periodontitis [7, 8], cutaneous leishmaniasis [9], obesity [10, 11], chronic kidney disease [12], Ulcerative colitis [13], and Crohn's disease [14, 15, 16]. The human microbiome appears to play a particularly important role in the development of Crohn's disease. Crohn's disease is an inflammatory bowel disease (IBD) with a prevalence of approximately 40 per 100,000 and 200 per 100,000 in children and adults, respectively [17].

Microbiomes present in the environment also serve important functions. The importance of ecological microbiomes is undisputed, due to their critical roles in fundamental processes such as nutrient cycling [18]. Due to differences in nutrient availability and environmental conditions, microbiomes in different environments are each characterized by unique structures and microbial compositions [19, 20, 21, 22]. Furthermore, each microbiome plays a unique role in the environment. For instance, the ocean microbiome generates half of the primary production on Earth [20]. The soil microbiome surrounding the root of plants has a great impact on plant fertility and growth [23]. Additionally, analyzing microbial communities in drinking water is one of the primary concerns of the environmental sciences [19].

The starting point in data collection for many of the above mentioned projects is 16S rRNA gene sequencing [24] of microbial samples, characterizing the taxonomic associations of the prokaryote or archaeal community members, which possess 16S genes. There are a number of features that make the 16S rRNA gene ideal for use as a taxonomic 'fingerprint' for microbiome composition characterization. First, the 16S rRNA gene is highly conserved across bacteria and archaea. Secondly, the gene consists of both conserved regions, for which universal species-independent PCR primers may be directed against, and nine hypervariable regions (V1-V9), along which species-specific sequence differences may accumulate to allow for differential identification of species [25]. Sequencing of specific hypervariable regions of the 16S rRNA gene may be performed for determination of microbial community composition [26]. After the 16S rRNA hypervariable regions are sequenced from a microbial sample, the obtained sequences are then processed using bioinformatics software (such as QIIME [27, 28], Mothur [29], or Usearch [30]) and clustered into groups of closely related sequences referred to as Operational Taxonomic Units (OTUs), which may then be used to assist in the functional profiling of microbial samples. Later in 1.2 we discuss the pros and cons of OTU features in details.

1.1 Machine learning for host/environment classification

Several recent studies predicted the environment or host phenotypes using 16S gene sequencing data for body-sites [31, 32], disease state [33, 34, 35, 36], ecological environment quality status prediction [37], and subject prediction for forensic science [38, 39]. In all, OTUs served as the main input feature for the down stream machine learning algorithms. Random Forest and then, ranking second, linear Support

Vector Machine (SVM) classifiers were reported as the most effective classification approaches in these studies [32, 35, 40].

Related prior work on body-site classification [31, 32] used the following datasets: Costello Body Habitat (CBH - 6 classes), Costello Skin Sites (CSS - 12 classes) [41], and Pei Body Site (PBS - 4 classes) [32]. An extensive comparison of classifiers for body-site classification over CBH, CSS, and PBS on top of OTU features has been performed by Statnikov et al [32]. The best accuracy levels measured by relative classifier information (RCI) achieved by using OTU features are reported as 0.784, 0.681, and 0.647 for CBH, CSS, and RCI respectively. Due to the insufficiency of the number of samples (on average 57 samples per class for CBH, CSS, and PBS) as well as the unavailability of raw sequences for some of the datasets mentioned above, instead of using the same dataset we replicate the state-of-the-art approach suggested in [32], i.e. Random Forest and SVM over OTU features for a larger dataset (Human Microbiome Project dataset). We then compare OTU features with k-mer representations. Working on a larger dataset allows for a more meaningful investigation and better training for deep learning approaches.

Detecting disease status based on 16S gene sequencing is becoming more and more popular, with applications in the prediction of Type 2 Diabetes [36] (patients: 53 samples, healthy:43 - Best accuracy: 0.67), Psoriasis (151 samples for 3 classes - Best accuracy: 0.225), IBD (patients: 49 samples, healthy:59 - Best AUC:0.95) [33], (patients: 91 samples, healthy: 58 samples - Best AUC:0.92) [34]. Similar to body-site classification datasets, the datasets used for disease prediction were also relatively small. In this paper, we use the Crohn's disease dataset [14] with 1359 samples for evaluating our proposed method and then compare it with the use of OTU features.

We focus on machine learning approaches for classification of environments or host phenotypes of 16S rRNA gene sequencing data, which is the most popular and cost-effective sequencing method for the characterization of microbiome to date [40]. Studies on the use of machine learning for predicting microbial phenotype instead of environments/host phenotype [42, 43], as well as predictions based on shotgun metagenomics and whole-genome microbial sequencing are beyond the scope of this paper, although we believe that one may easily adapt the proposed approach to shotgun metagenomics, similar to the study by Cui et al. on IBD prediction [44].

Recently, deep learning methods became popular in various applications of machine learning in bioinformatics [45, 46] and in particular in metagenomics [47]. However, to the best of our knowledge this paper is the first study exploring environment and host phenotype prediction from 16S rRNA gene sequencing data using deep learning approaches.

1.2 16S rRNA gene sequence representations

OTU representation

As reviewed in 1.1, prior machine learning works on environment/host phenotype prediction have been mainly using OTU representations as the input features to the learning algorithm. Almost all popular 16s rRNA gene sequence processing pipelines cluster sequences into OTUs based on their sequence similarities utilizing a variety of algorithms [28, 48]. QIIME allows OTU-picking using three different strategies: **(i) closed-reference OTU-picking:** sequences are compared against a marker gene database (e.g. Greengenes [49] or SILVA [50]) to be clustered into OTUs and then the sequences different from the reference genomes beyond a certain sequence identity threshold are discarded. **(ii) open-reference OTU-picking:** the remaining sequences after a closed-reference calling go through a de novo clustering. This allows for using the whole sequences as well as capturing sequences belonging to new communities which are absent in the reference databases [51]. **(iii) pure de novo OTU-picking:** sequences (or reads) are only compared among themselves and no reference database is used. The third strategy is more appropriate for novel species absent in the current reference. Although OTU clustering reduces the analysis of millions

of reads to working with only thousands of OTUs and simplifies the subsequent phylogeny estimation and multiple sequence alignment, OTU representations have several shortcomings: **(i)** All three OTU-picking strategies involve massive amounts of sequence alignments either to the reference genomes (in closed/opened-reference strategies) or to the sequences present in the sample (in open-reference and de novo strategies) which makes them very expensive [52] in comparison with reference-free/alignment-free representations. **(ii)** Overall sequence similarity is not a proper condition for grouping sequences and OTUs can be phylogenetically incoherent. For instance, a single mutation between two sequences is mostly ignored by OTU-picking algorithms. However, if the mutation does not occur within the sample, it might be a signal for assigning a new group. In addition, several mutations within a group most likely are not going to be tolerated by OTU-picking algorithms. However, having the same ratio across samples may suggest that the mutated sequences belong to the same group [48, 53]. **(iii)** The number of OTUs and even their contents are very sensitive to the pipeline and parameters, and this makes them difficult to reproduce [54].

k-mer representations

k-mer count vectors have been shown to be successful input features for performing machine learning algorithms on biological sequences for a variety of bioinformatics tasks [55]. In particular, recently k-mer count features have been largely used for taxonomic classifications of microbial metagenomics sequences [56, 57, 58, 59, 60, 61]. However, to the best of our knowledge, k-mer features have not been explored for phenotypical and environmental characterizations of 16S rRNA gene sequencing to simplify the classification pipeline. The advantages of using k-mer features over OTUs is later discussed in the discussion section.

In this paper we propose a new approach for environment/host phenotype prediction of 16S rRNA gene sequencing. Our approach is based on normalized k-mer distribution, which is fast, reference-free and alignment-free, while contributing in building accurate classifiers outperforming conventional OTU features in body-site identification and Crohn's disease classification tasks. We propose a bootstrapping framework to investigate the sufficiency of shallow sub-samples for the prediction of the phenotype of interest, which proves the sufficiency of short-length and shallow sequencing of 16S rRNA. In addition, we explore deep learning methods as well as classical approaches for the classification and show that in the presence of large datasets, deep learning can outperform classical methods. Furthermore, we explore PCA, t-SNE, and supervised deep representation learning for visualization of microbial samples/sequences of different phenotypes. We also show that k-mer features can be used to reliably predict representative 16S rRNA gene sequences belonging to 18 ecological environments and 5 organismal environments with high macro-F1s.

2 Material and Methods

2.1 Datasets

Body-site identification

We employ the metagenomic 16S rRNA gene sequence dataset provided by the NIH Human Microbiome Project (HMP) [62, 63]¹. In particular, we use processed, annotated 16S rRNA gene sequences of up to 300 healthy individuals, each sampled at 4 major body-sites (oral, airways, gut, vagina) and up to three time points. For each major body-site, a number of sub-sites were sampled. We focus on 5 body sub-sites: anterior nares (nasal) with 295 samples, saliva (oral) with 299 samples, stool (gut) with 325 samples, posterior fornix (urogenital) with 136 samples, and mid vagina (urogenital) with 137 samples, in

¹Available at <http://hmpdacc.org/HM16STR/>

total 1192 samples. These body-sites are selected to represent differing levels of spatial and biological proximity to one another, based on relevance to pertinent human health conditions potentially influenced by the human microbiome. In order to compare k-mer based approach with state-of-the-art OTU features we collect the closed-reference OTU representations of the same samples in HMP [63]² obtained using the QIIME pipeline [51].

Crohn's disease prediction

For the classification of Crohn's disease, we use the metagenomics 16S rRNA gene sequence dataset described in [14]³, which is currently the largest pediatric Crohn's disease dataset available. This dataset included annotated 16S rRNA gene sequence data for 731 pediatric (≤ 17 years old) patients with Crohn's disease and 628 samples verified as healthy or diagnosed with other diseases, making a total of 1359 samples. The 16S dataset was targeted towards the V4 hypervariable region of the 16S rRNA gene. Similar to the body-site dataset, in order to compare the k-mer based approach with the approach based on OTU features we collect the OTU representations of the same samples from Qiita repository⁴ obtained using QIIME pipeline [51].

Prediction of the environment for representative 16S rRNA gene sequences

MetaMetaDB provides a comprehensive dataset of representative 16S rRNA gene sequences of various ecological and organismal environments collected from existing 16S rRNA databases spanning almost 181 million raw sequences. In the MetaMetaDB pipeline, low-quality nucleotides, adaptors, ambiguous sequences, homopolymers, duplicates, and reads shorter than 200bp, as well as chimera have been removed and finally 16S rRNA sequences are clustered with 97% identity generating 1,241,213 representative 16S rRNA sequences marked by their environment [64]. MetaMetaDB divides its ecological environments into 34 categories and its organismal environments into 28 categories. We create three datasets which are subsets of MetaMetaDB to investigate the discriminative power of k-mers in predicting microbial habitability. Since the sequences in MetaMetaDB were already filtered and semi-identical sequences were removed, OTU-picking would not be relevant as it would result in an almost one-to-one mapping between the sequences and OTUs (we verified this using QIIME).

Ecological environment prediction: MetaMetaDB is imbalanced in terms of the number of representative sequences per environment. For this study, we pick the ecological environments with more than 10,000 samples, ending up with 18 classes of ecological environments: activated sludge, ant fungus garden, aquatic, bioreactor, bioreactor sludge, compost, food, food fermentation, freshwater, freshwater sediment, groundwater, hot springs, hydrocarbon, marine, marine sediment, rhizosphere, sediment, and soil⁵. We make two datasets out of the sequences in these environments: **ECO-18K** containing 1000 randomly selected instances per class (a total of 18K sequences) and **ECO-180K**, which is 10 times larger than **ECO-18K**, i.e. contains 10,000 randomly selected instances per class (a total of 180K sequences).

Organismal environment prediction: from the organismal environments in MetaMetaDB we select a subset containing gut microbiomes of 5 different organisms (bovine gut, chicken gut, human gut, mouse gut, termite gut) and down-sample each class to the size of the smallest class ending up having 620 sample per class (in total containing 3100 sequences) we call this dataset **5GUTS-3100**.

²Available at <https://qiita.ucsd.edu/study/descriptipn/1928>

³Available at: <https://www.ncbi.nlm.nih.gov/bioproject/PRJEB13679>

⁴Available at <https://qiita.ucsd.edu/study/description/1939>

⁵Datasets and descriptions are available at <http://mmdb.aori.u-tokyo.ac.jp/download.html>

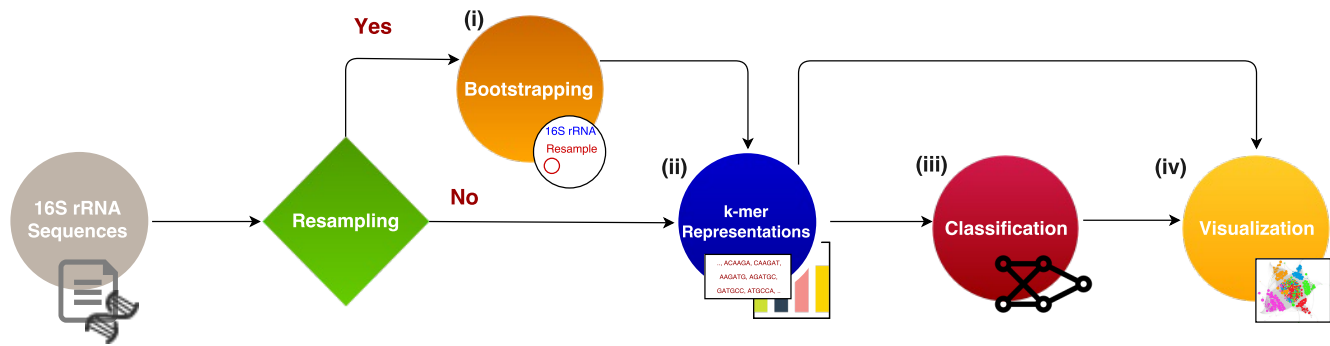


Figure 1. The components and the data flow in the MicroPheno computational workflow

2.2 MicroPheno computational workflow

We describe a computational workflow using deep learning and classical methods for classification of the environments/host phenotypes of microbial communities using k-mer frequency representations obtained from shallow sub-sampling of 16S rRNA gene sequences. We propose a bootstrapping framework to confirm the sufficiency of using a small portion of the sequences within a 16S rRNA sample for determining the underlying phenotype.

The MicroPheno computational workflow, as summarized in Figure 1, has the following steps: (i) to find the proper size for the sample such that it stays representative of the data and a stable k-mer profile, the 16S rRNA sequences go through a bootstrapping phase detailed in 2.2. (ii) Afterwards, the sub-sampled sequences will be used to produce k-mer representations of the samples. (iii) Then the k-mer representations will be used in classification algorithms, Deep Neural Networks (DNN), Random Forest (RF), and Linear SVM. (iv) Finally, the k-mer representations as well as the supervised representations trained using DNNs are used for visualization of the 16S rRNA gene sequences/samples. In what follows, these steps are explained in detail.

Bootstrapping

Confirming the sufficiency of only a small portion of 16S rRNA sequences for environment/phenotype classification is important because (i) sub-sampling reduces the preprocessing run-time, and (ii) more importantly, it proves that even a shallow 16S rRNA sequencing is enough. For this purpose we propose a resampling framework similar to bootstrapping to give us quantitative measures for finding the proper sampling size. Let $\theta_k(X_i)$ be the normalized k-mer distribution of X_i , a set of sequences in the i^{th} 16S rRNA sample. We investigate whether only a portion of X_i , which we represent as \tilde{x}_{ij} , i.e. j^{th} resample of X_i with sample size N , would be sufficient for producing a proper representation of X_i . To quantitatively find a sufficient sample size for X_i we propose the following criteria in a resampling scheme. **(i) Self-consistency:** resamples for a given size N from X_i produce consistent $\theta_k(\tilde{x}_{ij})$'s, i.e. resamples should have similar representations. **(ii) Representativeness:** resamples for a given size N from X_i produce $\theta_k(\tilde{x}_{ij})$'s similar to $\theta_k(X_i)$, i.e. similar to the case where all sequences are used.

We quantitatively define self-inconsistency and unrepresentativeness and seek parameter values that minimize them. We measure the **self-inconsistency** (\bar{D}_S) of the resamples' representations by calculating the average Kullback Leibler divergence among normalized k-mer distributions for N_R resamples (here $N_R=10$) with sequences of size N from the i^{th} 16S rRNA sample:

$$\bar{D}_{S_i}(N, k, N_R) = \frac{1}{N_R(N_R-1)} \sum_{\substack{p, q \in \{1, 2, \dots, N_R\} \\ (p \neq q)}} D_{KL}(\theta_k(\tilde{x}_{ip}), \theta_k(\tilde{x}_{iq})), \text{ where } |\tilde{x}_{il}| = N; \forall l \in \{1, 2, \dots, N_R\}.$$
 We calculate the average of the values of $\bar{D}_{S_i}(N, k, N_R)$ over the M different 16S rRNA samples:

$$\bar{D}_S(N, k, N_R) = \frac{1}{M} \sum_{i=1}^M \bar{D}_{S_i}(N, k, N_R).$$

We measure the **unrepresentativeness** (\bar{D}_R) of the resamples by calculating the average Kullback Leibler divergence between normalized k-mer distributions for N_R resamples ($N_R=10$) with size N and using all the sequences in X_i for the i^{th} 16S rRNA sample: $\bar{D}_{Ri}(N, k, N_R) = \frac{1}{N_R} \sum_{\forall p \in \{1, 2, \dots, N_R\}} D_{KL}(\theta_k(\tilde{x}_{ip}), \theta_k(X_i))$, where $|\tilde{x}_{ip}| = N; \forall i \in \{1, 2, \dots, N_R\}$. We calculate the average over $\bar{D}_{Ri}(N, k)$'s for the M 16S rRNA samples: $\bar{D}_R(N, k, N_R) = \frac{1}{M} \sum_{i=1}^M \bar{D}_{Ri}(N, k, N_R)$.

For the experiments on body-site and the dataset for Crohn's disease, we measure self-inconsistency \bar{D}_S and unrepresentativeness \bar{D}_R for $N_R = 10$ and $M = 10$ for any $8 \geq k \geq 3$ with sampling sizes ranging from 20 to 10000. Each point in Figure 5 represents the average of 100 ($M \times N_R$) resamples belonging to M randomly selected 16S rRNA samples, each of which is resampled $N_R = 10$ times. Since in the ecological and organismal datasets each sample is a single sequence, the bootstrapping step is skipped.

***k*-mer representation**

We propose l_1 normalized k-mer distribution of 16S rRNA gene sequences as input features for the downstream machine learning classification algorithms as well as visualization. Normalizing the representation allows for having a consistent representation even when the sampling size is changed. For each k-value we pick a sampling size that gives us a self-consistent and representative representation measured by $\bar{D}_S(N, k, N_R)$ and $\bar{D}_R(N, k, N_R)$ respectively as explained above. MicroPheno provides a parallel python implementation of k-mer distribution generation for a given sampling size.

Classification

Random Forests and linear SVM are the state-of-the-art classical approaches for categorical prediction on 16S rRNA sequencing [32, 35, 40] and in general for many machine learning problems in bioinformatics [65]. These two approaches, which are respectively instances of non-linear and linear classifiers, are both adopted in this study. In addition to these classical approaches, we also evaluate the performance of deep Neural Network classifiers in predicting environments and host phenotypes.

We evaluate and tune the model parameter in a stratified 10 fold cross-validation scheme. In order to ensure optimizing for both precision and recall we optimize the classifiers for the harmonic mean of precision and recall, i.e. F1. In particular, to give equal importance to the classification categories, specifically when we have imbalanced classes, we use macro-F1, which is the average of F1's over categories. Finally the evaluation metrics are averaged over the folds and the standard deviation is also reported.

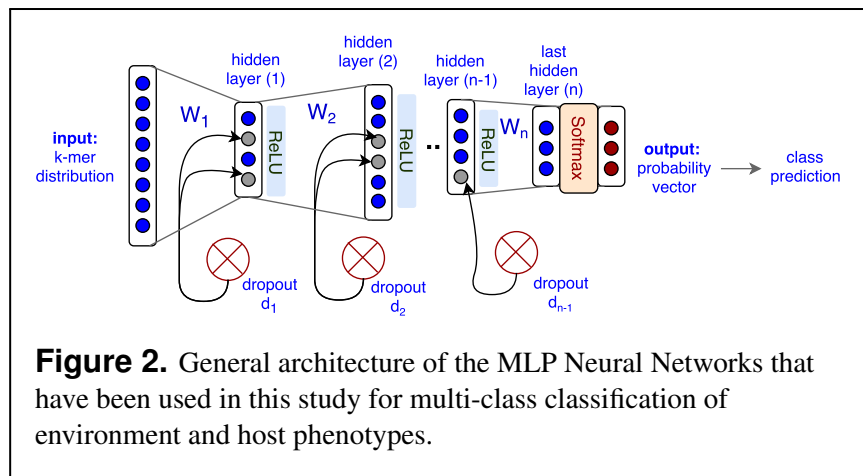
Classical learning algorithms: We use a one-versus-rest strategy for multi-class linear SVM [66] and tune parameter C , the penalty term for regularization. Random Forest [67] classifiers are tuned for (i) the number of decision trees in the ensemble, (ii) the number of features for computing the best node split, and (iii) the function to measure the quality of a split.

Deep learning: We use the Multi-Layer-Perceptrons (MLP) Neural Network architecture with several hidden layers using Rectified Linear Unit (ReLU) as the nonlinear activation function. We use softmax activation function at the last layer to produce the probability vector that can be regarded as representing posterior probabilities [68]. To avoid overfitting we perform early stopping and also use dropout at hidden layers [69]. A schematic visualization of our Neural Networks is depicted in Figure 2.

Our objective is minimizing the loss, i.e. cross entropy between output and the one-hot vector representation of the target class. The error (the distance between the output and the target) is used to update the network parameters via a Back-propagation algorithm using Adaptive Moment Estimation (Adam) as the optimizer [70].

We start with a single hidden layer and incrementally increase the number of layers with systematic exploration of the number of hidden units and dropout rates to find a proper architecture. We stop adding

layers when increasing the number of layers does not result in achieving a higher macro-F1 anymore. In addition, for the visualization of samples as explained later in 2.2 we use the output of the $(n - 1)^{th}$ hidden layer.



Implementations: MicroPheno uses implementations of Random Forest and SVM in the Python library scikit-learn [71] and deep Neural Networks are implemented in the Keras⁶ deep learning framework using TensorFlow back-end.

Visualization

In order to project 16S rRNA sequencing samples to 2D for visualization purposes, we explore Principal component analysis (PCA) [72] as well as t-Distributed Stochastic Neighbor Embedding (t-SNE) [73] as instances of respectively linear and non-linear dimensionality reduction methods. In addition, we explore the use of supervised deep representation learning in visualization of data [74], i.e. we visualize the activation function of the last hidden layer of the Neural Network trained for prediction of environments/host phenotypes to be compared with unsupervised methods. More details on these methods are provided as supplementary materials.

analysis (PCA) [72] as well as t-Distributed Stochastic Neighbor Embedding (t-SNE) [73] as instances of respectively linear and non-linear dimensionality reduction methods. In addition, we explore the use of supervised deep representation learning in visualization of data [74], i.e. we visualize the activation function of the last hidden layer of the Neural Network trained for prediction of environments/host phenotypes to be compared with unsupervised methods. More details on these methods are provided as supplementary materials.

3 Results

In this section, the results are organized based on datasets. As discussed in Section 2.2 we have several choices in each step in the computational workflow: choosing the value of k in k -mer, the sampling rate, and the classifiers. To explore the parameter space more systematically we followed the steps demonstrated in Figure 3. (i) In the first step for each value of $8 \geq k \geq 3$ we pick a stable sample size based on the output of bootstrapping. (ii) As the next step, we perform the classification task using tuned Random Forest for different k values and their selected sampling sizes based on bootstrapping. We selected Random Forest because we found it easy to tune in addition to the fact that it outperforms linear SVM in many cases. (iii) As the third step, for a selected k we investigate the role of sampling size (N) in classification. (iv) Finally, we compare different classifiers for the selected k and N . We also compare the performance of our proposed k -mer features with that of OTU features in classification tasks.

3.1 Body-site identification

- **(i) Bootstrapping for sampling rate selection for k -mers:** Higher k values require higher sampling rates to produce self-consistent and representative representations (Figure 5). For each k , the interval that \bar{D}_S and \bar{D}_R converge to their minimum values show a proper range for picking a sampling size resulting in self-consistent and representative representations.
- **(ii) Classification for different values of k with a sampling size selected based on the output of bootstrapping:** Interestingly, using only 3-mer features with a very low sampling rate ($\approx 20/15000=0.0013$) provides a relatively high performance for 5-way classification. The value of

⁶<https://keras.io/>

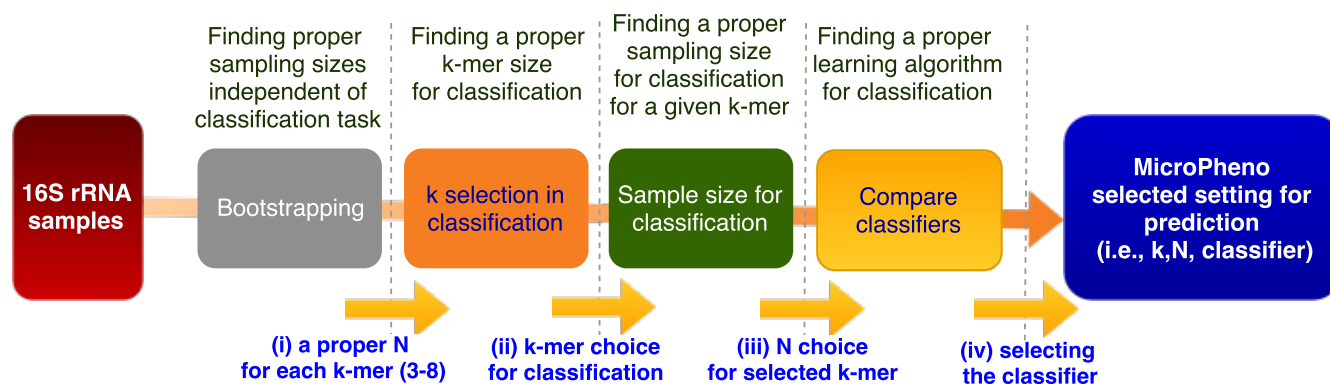


Figure 3. Steps we take to explore parameters for the representations and how we choose the classifier for prediction of the phenotype of interest in this study

macro-F1 increases with the value of k from 3 to 6, but increasing k further than that does not have any additional effect on macro-F1 (Table 1, step (ii)).

- **(iii) Exploring the sampling size (N) for a selected k-mer:** For a selected k -value ($k=6$), using Random Forest classifier for different sampling sizes are presented in Table 1, step (iii). The results suggest that changing the sampling size from 0.6% to 100% of the sequences will not change the classification results significantly, suggesting that in body-site identification, a very shallow sub-sampling of the sequences is sufficient for a reliable prediction. Using more sequences does not necessarily increase the discriminative power and may even result in over-fitting. We selected a sampling size of 5000 for 6-mers for comparison between classifiers in the next step.
- **(iv) Comparison of classifiers for the selected N, k:** For selected values of k , N , the results of the body-site prediction task using Random Forest, SVM, and Neural Network classifiers are provided (Table 1, step (iv)). Random Forest classifier obtained the top macro-F1 (0.84) for this 5-way classification.

The confusion matrix in Figure 4 shows that the most difficult decision for the classifier is to distinguish between mid vagina and posterior fornix, both of which are urogenital body-sites. The visualizations of body-site samples obtained through using PCA, t-SNE, and t-SNE on the activation function of the last layer of the trained 5-layered Neural Network are presented in Figure 6. These results suggest that supervised training of representations using Neural Networks provides a non-linear transformation of data that can discriminate between dissimilar body-sites with reasonably accuracy. As shown in the last row of Table 1, combining the urogenital body-sites increases the macro-F1 to 0.99 ± 0.01 using the Neural Network.

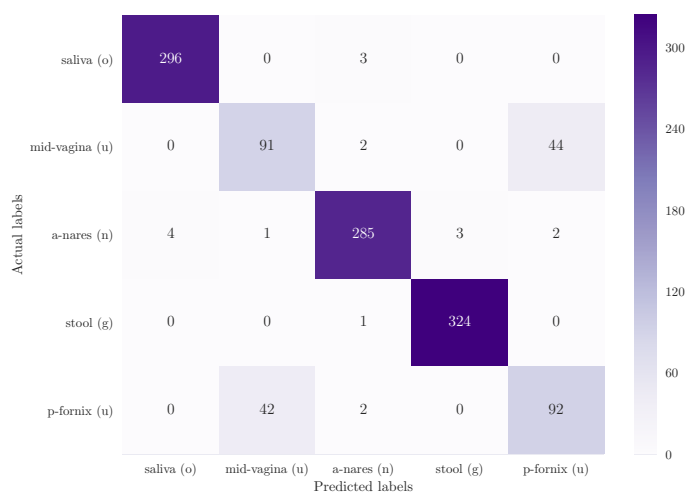


Figure 4. The confusion matrix for the classification of 5 major body-sites, using Random Forest classifier in a 10-fold cross-validation scheme. The presented body-sites are saliva (**o**: oral), mid-vagina (**u**: urogenital), anterior nares (**n**: nasal), stool (**g**: gut), and posterior fornix (**u**: urogenital).

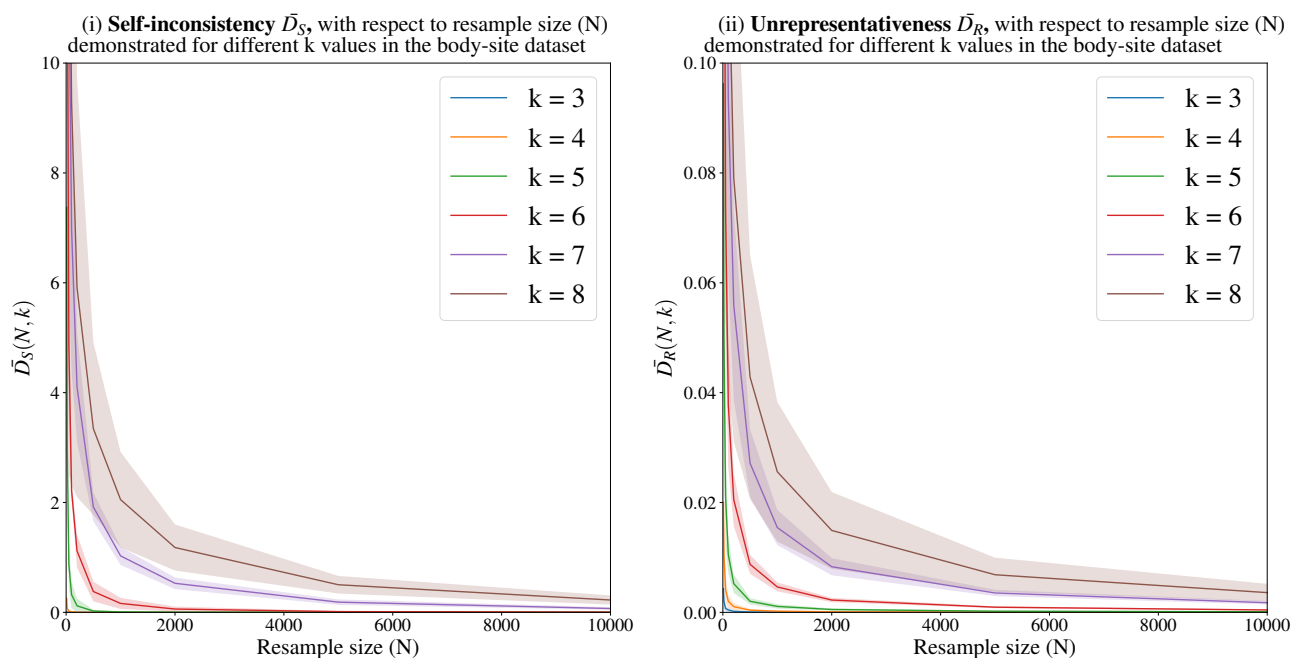


Figure 5. Measuring (i) self-inconsistency (\bar{D}_S), and unrepresentativeness (\bar{D}_R) for the body-site dataset. Each point represents an average of 100 resamples belonging to 10 randomly selected 16S rRNA samples. Higher k values require higher sampling rates to produce self-consistent and representative samples.

Step	Representation	Resample size/ \approx 15000/fasta	Classifier	Micro-metrics (averaged over samples)			Macro-metrics (averaged over classes)		
				Precision	Recall	F1	Precision	Recall	F1
(ii)	3-mers	20	RF	0.84 \pm 0.02	0.84 \pm 0.02	0.84 \pm 0.02	0.75 \pm 0.03	0.75 \pm 0.03	0.74 \pm 0.03
	4-mers	100		0.86 \pm 0.03	0.86 \pm 0.03	0.86 \pm 0.03	0.77 \pm 0.03	0.77 \pm 0.03	0.77 \pm 0.03
	5-mers	500		0.89 \pm 0.02	0.89 \pm 0.02	0.89 \pm 0.02	0.82 \pm 0.03	0.82 \pm 0.03	0.82 \pm 0.03
	6-mers	2000		0.91 \pm 0.03	0.91 \pm 0.03	0.91 \pm 0.03	0.85 \pm 0.05	0.85 \pm 0.04	0.84 \pm 0.05
	7-mers	5000		0.91 \pm 0.03	0.91 \pm 0.03	0.91 \pm 0.03	0.85 \pm 0.05	0.85 \pm 0.05	0.85 \pm 0.05
	8-mers	8000		0.9 \pm 0.03	0.9 \pm 0.03	0.9 \pm 0.03	0.85 \pm 0.05	0.84 \pm 0.05	0.84 \pm 0.05
(iii)	6-mers	100	RF	0.89 \pm 0.02	0.89 \pm 0.02	0.89 \pm 0.02	0.82 \pm 0.04	0.82 \pm 0.03	0.81 \pm 0.03
		1000		0.9 \pm 0.03	0.9 \pm 0.03	0.9 \pm 0.03	0.83 \pm 0.04	0.83 \pm 0.04	0.83 \pm 0.04
		2000		0.91 \pm 0.03	0.91 \pm 0.03	0.91 \pm 0.03	0.85 \pm 0.05	0.85 \pm 0.04	0.84 \pm 0.05
		5000		0.9 \pm 0.03	0.9 \pm 0.03	0.9 \pm 0.03	0.85 \pm 0.04	0.84 \pm 0.04	0.84 \pm 0.04
		10000		0.9 \pm 0.03	0.9 \pm 0.03	0.9 \pm 0.03	0.84 \pm 0.05	0.84 \pm 0.05	0.84 \pm 0.05
		All sequences		0.9 \pm 0.03	0.9 \pm 0.03	0.9 \pm 0.03	0.84 \pm 0.05	0.84 \pm 0.04	0.84 \pm 0.05
(iv)	6-mers	5000	RF	0.9 \pm 0.03	0.9 \pm 0.03	0.9 \pm 0.03	0.85 \pm 0.04	0.84 \pm 0.04	0.84 \pm 0.04
			SVM	0.86 \pm 0.02	0.86 \pm 0.02	0.86 \pm 0.02	0.76 \pm 0.06	0.76 \pm 0.03	0.74 \pm 0.04
			DNN-5	0.87 \pm 0.01	0.87 \pm 0.01	0.87 \pm 0.01	0.79 \pm 0.02	0.79 \pm 0.03	0.79 \pm 0.02
-	6-mers	5000	DNN-4 (4 classes)	0.99 \pm 0.01	0.99 \pm 0.01	0.99 \pm 0.01	0.99 \pm 0.01	0.99 \pm 0.01	0.99 \pm 0.01

Table 1. The results for classification of major body-sites using k-mer representations. The set of rows matches the steps (ii to iv) mentioned in Figure 3, i.e k-mer selection, N (sample size) selection, and finally selection of the classifier. The classifiers (Random Forest, Support Vector Machine and Neural Network classifiers) are tuned and evaluated in a stratified 10x fold cross-validation setting. The last row shows the Neural Network's performance in the classification of body-sites when the urogenital body-sites are combined.

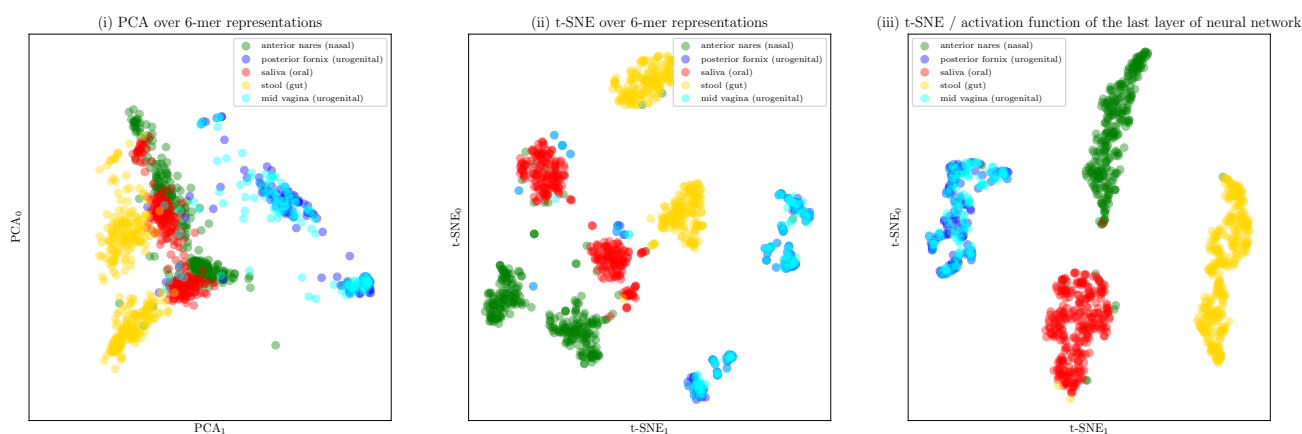


Figure 6. Visualization of the body-site dataset using different projection methods: (i) PCA over 6-mer distributions with unsupervised training, (ii) t-SNE over 6-mer distributions with unsupervised training, (iii) visualization of the activation function of the last layer of the trained Neural Network (projected to 2D using t-SNE).

Features	Classifiers	Micro-metrics (averaged over samples)			Macro-metrics (averaged over classes)		
		Precision	Recall	F1	Precision	Recall	F1
6-mer features (size: 4096)	RF	0.9 ± 0.03	0.9 ± 0.03	0.9 ± 0.03	0.85 ± 0.04	0.84 ± 0.04	0.84 ± 0.04
	SVM	0.86 ± 0.02	0.86 ± 0.02	0.86 ± 0.02	0.76 ± 0.06	0.76 ± 0.03	0.74 ± 0.04
OTU features (size: 20589)	RF	0.89 ± 0.02	0.89 ± 0.02	0.89 ± 0.02	0.83 ± 0.03	0.83 ± 0.03	0.83 ± 0.03
	SVM	0.85 ± 0.03	0.85 ± 0.03	0.85 ± 0.03	0.77 ± 0.05	0.78 ± 0.04	0.76 ± 0.04

Table 2. Comparison of k-mer and OTU features in body-site identification.

Comparison of k-mer and OTU features in body-site identification

A comparison between OTU features and k-mer representations in body-site identification is presented in Table 2. For this comparison, Random Forest classifier (as an instance of non-linear classifier) and linear SVM (as an instance of linear classifier) have been used, tuned and evaluated in a stratified 10x fold cross-validation setting. Our results suggest that for both k-mer features and OTUs, Random Forest is the best choice. In addition, with almost $\frac{1}{5}$ of the size of OTU features and in spite of being considerably less expensive to calculate, k-mer marginally outperforms OTU features in body-site identification.

3.2 Crohn's disease prediction

- **(i) Bootstrapping for sampling rate selection for k-mers:** Similar to bootstrapping for the body-site dataset (Figure 5), \bar{D}_S and \bar{D}_R for different values of k with respect to sampling sizes were calculated. As the structure of the curve is similar to the body-site dataset, in order to avoid redundancy the figure is provided as supplementary material.
- **(ii) Classification for different values of k with a sampling size selected based on the output of bootstrapping:** Choosing k=6 with a sampling size of 2000 ($\approx 2,000/38,000 = 0.05$) provides us a macro-F1 of 0.75 which is the minimum k with top performance (Table 3, step (ii)).
- **(iii) Exploring the sampling size (N) for a selected k-mer:** For a selected k-value (k=6), using Random Forest classifier for different sampling sizes are presented in Table 3, step (iii). Increasing the sampling size from 100 ($100/38000=0.003$) to 5000 ($5000/38000=0.13$) increases the macro-F1 from 0.7 to 0.75. In addition, using all sequences instead of 0.13 of them in each sample, does not increase the discriminative power.
- **(iv) Comparison of classifiers for the selected N, k:** For selected values of k, N, the results of the Crohn's disease prediction using Random Forest, SVM, and Neural Network classifiers are

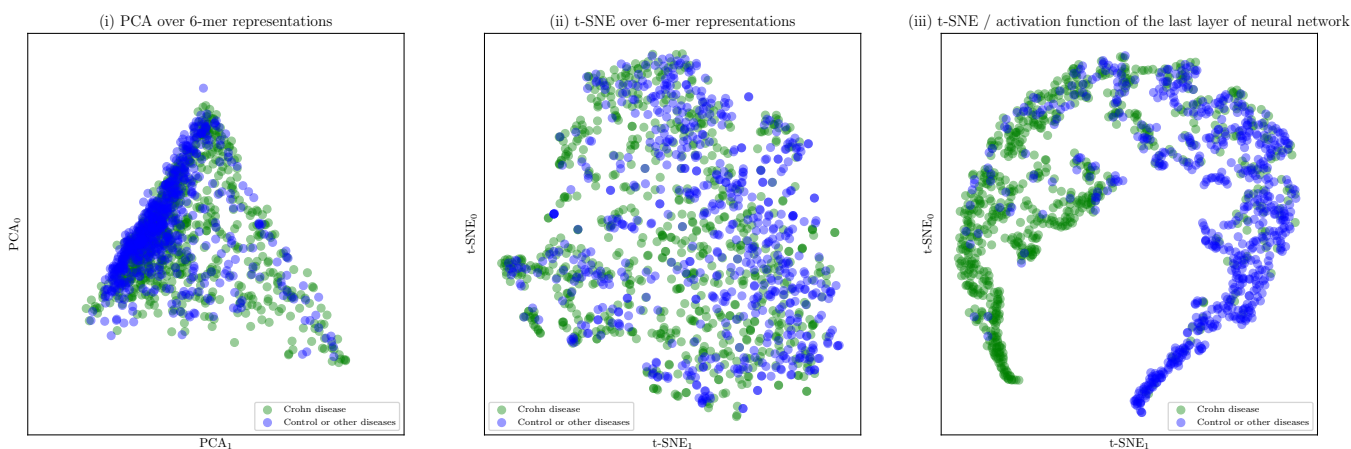


Figure 7. Visualization of the Crohn's disease dataset using different projection methods: (i) PCA over 6-mer distributions with unsupervised training, (ii) t-SNE over 6-mer distributions with unsupervised training, (iii) visualization of the activation function of the last layer of the trained Neural Network (projected to 2D using t-SNE).

Step	Representation	Resample size/ \approx 38000/fasta	Classifier	Micro-metrics (averaged over samples)			Macro-metrics (averaged over classes)		
				Precision	Recall	F1	Precision	Recall	F1
(ii)	3-mers	20	RF	0.62 \pm 0.05	0.62 \pm 0.05	0.62 \pm 0.05	0.62 \pm 0.05	0.61 \pm 0.05	0.61 \pm 0.05
	4-mers	100		0.7 \pm 0.05	0.7 \pm 0.05	0.7 \pm 0.05	0.69 \pm 0.05	0.69 \pm 0.05	0.69 \pm 0.05
	5-mers	500		0.74 \pm 0.05	0.74 \pm 0.05	0.74 \pm 0.05	0.74 \pm 0.05	0.74 \pm 0.05	0.74 \pm 0.05
	6-mers	2000		0.75 \pm 0.04	0.75 \pm 0.04	0.75 \pm 0.04	0.75 \pm 0.04	0.75 \pm 0.04	0.75 \pm 0.04
	7-mers	5000		0.75 \pm 0.04	0.75 \pm 0.04	0.75 \pm 0.04	0.75 \pm 0.04	0.75 \pm 0.04	0.75 \pm 0.04
	8-mers	8000		0.75 \pm 0.04	0.75 \pm 0.04	0.75 \pm 0.04	0.75 \pm 0.04	0.75 \pm 0.04	0.75 \pm 0.04
(iii)	6-mers	100	RF	0.71 \pm 0.04	0.71 \pm 0.04	0.71 \pm 0.04	0.71 \pm 0.04	0.7 \pm 0.04	0.7 \pm 0.04
		1000		0.75 \pm 0.04	0.75 \pm 0.04	0.75 \pm 0.04	0.76 \pm 0.04	0.75 \pm 0.04	0.75 \pm 0.04
		2000		0.75 \pm 0.04	0.75 \pm 0.04	0.75 \pm 0.04	0.75 \pm 0.04	0.75 \pm 0.04	0.75 \pm 0.04
		5000		0.76 \pm 0.04	0.76 \pm 0.04	0.76 \pm 0.04	0.76 \pm 0.04	0.75 \pm 0.04	0.75 \pm 0.04
		10000		0.76 \pm 0.04	0.76 \pm 0.04	0.76 \pm 0.04	0.76 \pm 0.05	0.75 \pm 0.04	0.75 \pm 0.05
		All sequences		0.76 \pm 0.04	0.76 \pm 0.04	0.76 \pm 0.04	0.76 \pm 0.05	0.75 \pm 0.04	0.75 \pm 0.05
(iv)	6-mers	5000	RF	0.76 \pm 0.04	0.76 \pm 0.04	0.76 \pm 0.04	0.76 \pm 0.04	0.75 \pm 0.04	0.75 \pm 0.04
			SVM	0.68 \pm 0.04	0.68 \pm 0.04	0.68 \pm 0.04	0.68 \pm 0.04	0.67 \pm 0.04	0.67 \pm 0.04
			DNN-7	0.7 \pm 0.02	0.7 \pm 0.02	0.7 \pm 0.02	0.7 \pm 0.03	0.7 \pm 0.02	0.7 \pm 0.03

Table 3. The results for classification of Crohn's disease prediction using k-mer representations. The set of rows matches the steps (ii to iv) mentioned in Figure 3, i.e k-mer selection, N (sample size) selection, and finally selection of the classifier. The classifiers (Random Forest, Support Vector Machine and Neural Network classifiers) are tuned and evaluated in a stratified 10xfold cross-validation setting.

provided (Table 3, step (iv)). Random Forest classifier obtained the top macro-F1 (0.75) for this binary classification.

The projection of the Crohn's disease dataset using with PCA and t-SNE over raw k-mer representations, as well as t-SNE over the representation learned though the supervised learning of the Neural Network are visualized in Figure 7. The trained Neural Network provides a non-linear transformation of data that can identify subjects diagnosed with Crohn's disease with reasonable performance.

Comparison of k-mer and OTU features in Crohn's disease classification

For a comparison between OTU features and our proposed k-mer representations in detecting Crohn's disease from metagenomic samples, the Random Forest classifier (as an instance of non-linear classifiers) and linear SVM (as an instance of linear classifiers) were tuned and evaluated in a stratified 10xfold cross-validation. For both k-mer features and OTUs, the Random Forest classifier performed best (Table 4). In addition, even though only half of the number of features were used and in spite of being considerably

Features	Classifiers	Micro-metrics (averaged over samples)			Macro-metrics (averaged over classes)		
		Precision	Recall	F1	Precision	Recall	F1
6-mer features (size: 4096)	RF	0.76 ± 0.04	0.76 ± 0.04	0.76 ± 0.04	0.76 ± 0.04	0.75 ± 0.04	0.75 ± 0.04
	SVM	0.68 ± 0.04	0.68 ± 0.04	0.68 ± 0.04	0.68 ± 0.04	0.67 ± 0.04	0.67 ± 0.04
OTU features (size: 9511)	RF	0.74 ± 0.04	0.74 ± 0.04	0.74 ± 0.04	0.74 ± 0.04	0.74 ± 0.04	0.74 ± 0.04
	SVM	0.68 ± 0.04	0.68 ± 0.04	0.68 ± 0.04	0.68 ± 0.04	0.68 ± 0.04	0.68 ± 0.04

Table 4. Comparison of k-mers and OTU features in the detection of the Crohn’s disease phenotype. For this comparison, Random Forest classifier (as an instance of non-linear classifiers) and linear SVM (as an instance of linear classifiers) have been used. The classifiers are tuned and evaluated in a stratified 10xfold cross-validation setting.

Step	Representation	Dataset	Classifier	Micro-metrics (averaged over samples)			Macro-metrics (averaged over classes)		
				Precision	Recall	F1	Precision	Recall	F1
(ii)	3-mers	ECO-18K	RF	0.6 ± 0.01	0.6 ± 0.01	0.6 ± 0.01	0.63 ± 0.02	0.6 ± 0.01	0.57 ± 0.01
	4-mers			0.67 ± 0.01	0.67 ± 0.01	0.67 ± 0.01	0.7 ± 0.01	0.67 ± 0.01	0.65 ± 0.01
	5-mers			0.72 ± 0.01	0.72 ± 0.01	0.72 ± 0.01	0.74 ± 0.01	0.72 ± 0.01	0.71 ± 0.01
	6-mers			0.75 ± 0.01	0.75 ± 0.01	0.75 ± 0.01	0.76 ± 0.01	0.75 ± 0.01	0.73 ± 0.01
	7-mers			0.74 ± 0.01	0.74 ± 0.01	0.74 ± 0.01	0.76 ± 0.01	0.74 ± 0.01	0.73 ± 0.01
	8-mers			0.72 ± 0.01	0.72 ± 0.01	0.72 ± 0.01	0.74 ± 0.01	0.72 ± 0.01	0.71 ± 0.01
(iv)	6-mers	ECO-18K	RF	0.75 ± 0.01	0.75 ± 0.01	0.75 ± 0.01	0.76 ± 0.01	0.75 ± 0.01	0.73 ± 0.01
			SVM	0.79 ± 0.01	0.79 ± 0.01	0.79 ± 0.01	0.79 ± 0.01	0.79 ± 0.01	0.79 ± 0.01
			DNN-3	0.78 ± 0.01	0.78 ± 0.01	0.78 ± 0.01	0.78 ± 0.01	0.78 ± 0.01	0.78 ± 0.01
(iv)	6-mers	ECO-180K (10x larger)	RF	0.83 ± 0.0	0.83 ± 0.0	0.83 ± 0.0	0.84 ± 0.0	0.83 ± 0.0	0.83 ± 0.0
			SVM	0.86 ± 0.0	0.86 ± 0.0	0.86 ± 0.0	0.87 ± 0.01	0.86 ± 0.0	0.86 ± 0.0
			DNN-5	0.88 ± 0.0	0.88 ± 0.0	0.88 ± 0.0	0.88 ± 0.0	0.88 ± 0.0	0.88 ± 0.0

Table 5. The results for the task of selecting between 18 ecological environments. The classifiers (Random Forest, Support Vector Machine and Neural Network classifiers) are tuned and evaluated in a stratified 10xfold cross-validation setting over both ECO-18K and ECO-180K datasets.

less expensive to calculate, k-mers marginally outperformed OTU features in the detection of Crohn’s dis

3.3 Ecological environment prediction

- **(ii) Classification for different values of k:** As stated before, for the ecological and organismal datasets we do not need to perform resampling as we classify single 16S rRNA gene representative sequences. We thus can skip steps (i) and (iii) (Figure 3). Step (ii) in Table 5, shows the effect of k in the performance of the classification of the 18 environments for the ECO-18K dataset. The result shows that k=6 provides a better classification performance with a macro-F1 of 0.73 which is relatively high for a 18-way classification (has a mere 0.06 chance of randomly occurring).
- **(iv) Comparison of classifiers for the selected k:** For selected values of k the results of the environment prediction using Random Forest, SVM, and Neural Network classifiers are provided (Table 5, step (iv), ECO-18K dataset). SVM classifier obtained the top macro-F1 (0.79) for 18-way classification. In order to see the effect of increasing the number of data points in classification performance we repeat the classifier comparison (step iv) for the ECO-180K dataset. The results are summarized in Table 5 showing that feeding more training instances results in better training for the deep learning approach, outperforming SVM and achieving a macro-F1 of 0.88, which is very high for a 18-way classification framework.

In training the Neural Networks for the ECO-18K dataset, increasing the number of hidden layers from 3 to more did not help result in improvements. However, using the ECO-180K dataset, which is 10 times larger, allowed us to train a deeper network increasing the macro-F1 by 5 percent going from 3 layers to 5 layers. Increasing the number of layers further did not result in any improvements.

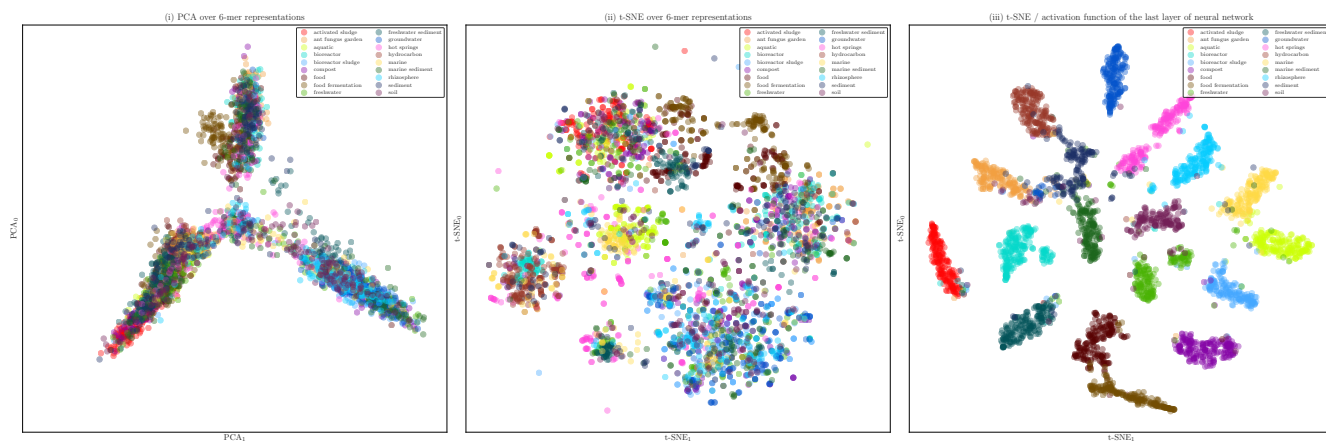


Figure 8. Visualization of 18 ecological microbial environments using different projection methods: (i) PCA over 6-mer distributions with unsupervised training, (ii) t-SNE over 6-mer distributions with unsupervised training, (iii) visualization of the activation function of the last layer of the trained Neural Network (projected to 2D using t-SNE).

Neural Network Visualization

Visualizations of representative 16S rRNA gene sequences in 18 ecological environments obtained through using PCA, t-SNE, and t-SNE on the activation function of the last layer of the trained Neural Network are presented in 8. For ease of visualization, we randomly picked 100 samples per class. These results suggest that supervised training of representations using Neural Networks provides a non-linear transformation of data containing information about high-level similarities between environments in the sub-plot on the right (scatter plot (iii)), where such structures appeared in the visualization only when more hidden layers were used: **(i)** On the left, the environments containing water are clustered in a dense neighborhood: marine, aquatic, freshwater, hot springs, bioreactor sludge (description in the dataset: The sludge inside the bioreactor that treats wastewater.), groundwater, and rhizosphere (an environment where plants, soil, water, microorganisms, and nutrients meet and interact). **(ii)** In the middle we have environments labeled related to sediment: sediment, freshwater sediment, marine sediment, and soil and. **(iii)** Ant fungus garden is close to hydrocarbon. A previous study confirms that chemical analysis of ant fungus gardens reports the presence of hydrocarbons, a class of chemical compounds found commonly on insect cuticles [75]. **(iv)** Environments containing food like food, food fermentation, and compost are at the bottom of the plot. **(i)** Finally, artificial and industrial environments like bioreactor and activated sludge are clustered on the left of the sub-plot.

3.4 Organismal environment identification

- **(ii) Classification for different values of k:** The results show that $k=6$ and 7 provide a high classification macro-F1 of 0.87 for 5 classes (0.2 chance of randomly occurring), step (ii) in Table 6.
- **(iv) Comparison of classifiers for the selected k:** For selected values of k the results of the organism prediction using Random Forest, SVM, and Neural Network classifiers are provided (Table 6, step (iv)). SVM classifier obtained the top macro-F1 (0.88) for 5-way classification.

4 Discussion

In this work, we present MicroPheno, a new approach for environments and host phenotypes prediction using normalized k -mer distribution of 16S rRNA gene sequences over shallow sub-samples. We divide our discussion of the results of this study into three main components: (1) the use of k -mers versus OTUs,

Step	Representation	Dataset	Classifier	Micro-metrics (averaged over samples)			Macro-metrics (averaged over classes)		
				Precision	Recall	F1	Precision	Recall	F1
(ii)	3-mers	5GUTS-3100	RF	0.8 ± 0.02	0.8 ± 0.02	0.8 ± 0.02	0.8 ± 0.02	0.8 ± 0.02	0.79 ± 0.02
	4-mers			0.84 ± 0.01	0.84 ± 0.01	0.84 ± 0.01	0.84 ± 0.01	0.84 ± 0.01	0.83 ± 0.01
	5-mers			0.86 ± 0.02	0.86 ± 0.02	0.86 ± 0.02	0.86 ± 0.02	0.86 ± 0.02	0.85 ± 0.02
	6-mers			0.87 ± 0.01	0.87 ± 0.01	0.87 ± 0.01	0.87 ± 0.01	0.87 ± 0.01	0.87 ± 0.01
	7-mers			0.87 ± 0.01	0.87 ± 0.01	0.87 ± 0.01	0.88 ± 0.02	0.87 ± 0.01	0.87 ± 0.01
	8-mers			0.86 ± 0.01	0.86 ± 0.01	0.86 ± 0.01	0.87 ± 0.01	0.86 ± 0.01	0.86 ± 0.01
(iv)	6-mers	5GUTS-3100	RF	0.87 ± 0.01	0.87 ± 0.01	0.87 ± 0.01	0.87 ± 0.01	0.87 ± 0.01	0.87 ± 0.01
			SVM	0.88 ± 0.02	0.88 ± 0.02	0.88 ± 0.02	0.89 ± 0.01	0.88 ± 0.02	0.88 ± 0.02
			DNN-5	0.87 ± 0.01	0.87 ± 0.01	0.87 ± 0.01	0.87 ± 0.01	0.87 ± 0.01	0.87 ± 0.01

Table 6. The results for the task of classifying 5 organismal environments belonging to 5 organisms' gut. The classifiers (Random Forest, Support Vector Machine and Neural Network classifiers) are tuned and evaluated in a stratified 10x fold cross-validation setting over 5GUTS-3100 dataset.

(2) the benefits of shallow sub-sampling, and (3) classical methods versus the deep learning approach.

4.1 K-mers versus OTUs

In order to evaluate MicroPheno, we compared our proposed k-mer representations with OTU features in two tasks of body-site identification and Crohn's disease classification. We replicated the state-of-the-art approach, i.e. Random Forest over OTU features, on datasets larger than those that had been previously explored. We showed that k-mer features outperform conventional OTUs, while having several advantages over OTUs, as listed below:

- k-mer representations are easy to compute at no computational cost for any type of alignment to references or tasks of finding pair-wise sequence similarity within samples as needed in OTU-picking pipelines. Just to get an idea of the computational efficiency of k-mers calculation in comparison with OTUs, note that OTU picking for the Crohn's disease dataset of 1359 samples takes more than 5 hours using 5 threads, while 6-mer distribution calculation took around 5 minutes using the optimal sampling size (N) for classification using the same number of threads.
- Taxonomy-independent analysis is often the preferred approach for amplicon sequencing when the samples contains unknown taxa. k-mer features can be used without making any assumptions about the taxonomy. However, OTU-picking pipelines make assumptions about the taxonomy as discussed in 1.2; therefore they can even be phylogenetically incoherent.
- k-mer distribution is a well-defined and stable representation, while OTUs are sensitive to the pipeline and the parameters
- Sequence similarities are naturally incorporated in the k-mer representations for the downstream learning algorithm, but with grouping sequences into certain categories, the sequence similarities between OTUs are ignored
- No need for full length sequencing of 16S rRNA: using short k-mers (k=6) suggests that even cheaper technologies offering very short length sequencing of 16S can be sufficient for the prediction of the phenotype of interest.

NO-FREE-LUNCH; and the main disadvantage of k-mer features over OTUs is that, using short k-mers make it more difficult to trace the relevant taxa to the phenotype of interest. When such an analysis is needed using OTUs or increasing the size of k would be an alternative solution. However, as long as prediction is concerned, the k-mer representation seems to be the best choice for an accurate and rapid detection/diagnosis over 16S rRNA sequencing samples.

4.2 Shallow sub-sampling

We proposed a bootstrapping framework to investigate the consistency and representativeness of k-mer distribution for different sampling rates. Our results suggest that depending on the k-mer size even very low sub-sampling rates (0.001 to 0.1) (for k between 3 to 7) not only can provide a consistent representation, but can also result in better predictions while possibly avoiding overfitting. Setting aside the save in preprocessing time as a natural benefit of sampling, this result also suggest that at least for similar phenotypes of interest, shallow sequencing of the microbial community would be sufficient for an accurate prediction.

4.3 Classical classifiers versus Deep learning

To the best of our knowledge this paper for the first time explores the use of deep learning for environment/host phenotype prediction of 16S rRNA sequences. Studying the role of dataset size in the classification of ecological environments showed that for large datasets using deep learning provides us with more accurate predictions. However, when the number of samples are not large enough, Random Forests performs better on both OTUs and k-mer features. In addition, we observed that in the case of classification over representative sequences as opposed to samples (pool of sequences), SVM works better than Random Forest classifier.

Another advantage of using deep learning in classification was that supervised training of a proper representation of data results in a more discriminative representation for the downstream visualization compared to the unsupervised methods (PCA and t-SNE on the raw k-mer distributions). In the cases of body-site identification and more clearly in ecological environment classification, the model was able to extract more high-level similarities between the environments as detailed in 3.3.

5 Conclusion

A new approach for environment/host phenotype prediction on 16S rRNA gene sequencing has been presented based on k-mer representations of shallow sub-samples, outperforming the computationally costly OTU features in the two tasks of body-site identification and Crohn's disease classification. This result also suggests the sufficiency of a shallow and short length sequencing of 16S rRNA sequences for phenotype prediction purposes. Deep learning methods as well as classical approaches were explored for the classification of the the environments/phenotypes. In addition, we showed that k-mer features can reliably predict representative 16S rRNA gene sequences of 18 ecological environments, and 5 organismal environments with high macro-F1 scores of 0.88 and 0.87. We showed that in the presence of large datasets, deep learning can outperform classical methods such as Random Forest and SVM. Furthermore, PCA, t-SNE, and supervised deep representation learning were explored in this paper for visualization of microbial samples/sequences of different phenotypes.

References

1. Lynch, S. V. & Pedersen, O. The Human Intestinal Microbiome in Health and Disease. *New Engl. J. Medicine* **375**, 2369–2379 (2016). URL <http://www.nejm.org/doi/10.1056/NEJMra1600266>. DOI 10.1056/NEJMra1600266.
2. Marsland, B. J., Yadava, K. & Nicod, L. P. The airway microbiome and disease. *Chest* **144**, 632–637 (2013). DOI 10.1378/chest.12-2854.
3. Arrieta, M.-C. *et al.* Early infancy microbial and metabolic alterations affect risk of childhood asthma. *Sci. Transl. Medicine* **7**, 307ra152–307ra152 (2015). URL <http://stm.sciencemag.org/scitransmed/7/307/307ra152.full.pdf>. DOI 10.1126/scitranslmed.aab2271.
4. Saulnier, D. M. *et al.* Gastrointestinal microbiome signatures of pediatric patients with irritable bowel syndrome. *Gastroenterol.* **141**, 1782–1791 (2011). DOI 10.1053/j.gastro.2011.06.072.

5. Cho, I. & Blaser, M. J. The human microbiome: At the interface of health and disease (2012). DOI 10.1038/nrg3182. [NIHMS150003](#).
6. Cammarota, G., Ianiro, G. & Gasbarrini, A. Fecal microbiota transplantation for the treatment of clostridium difficile infection: A systematic review (2014). DOI 10.1097/MCG.000000000000046.
7. Deng, Z. L., Szafranski, S. P., Jarek, M., Bhujju, S. & Wagner-Döbler, I. Dysbiosis in chronic periodontitis: Key microbial players and interactions with the human host. *Sci. Reports* **7**, 1–13 (2017). DOI 10.1038/s41598-017-03804-8.
8. Jorth, P. *et al.* Metatranscriptomics of the Human Oral Microbiome during Health and Disease. *mBio* **5**, e01012–14–e01012–14 (2014). URL <http://mbio.asm.org/content/5/2/e01012-14.full.pdf+html>{%}5Cnpapers2://publication/doi/10.1128/mBio.01012-14. DOI 10.1128/mBio.01012-14.
9. Gimblet, C. *et al.* Cutaneous Leishmaniasis Induces a Transmissible Dysbiotic Skin Microbiota that Promotes Skin Inflammation. *Cell Host Microbe* **22**, 13–24.e4 (2017). DOI 10.1016/j.chom.2017.06.006.
10. Turnbaugh, P. J., Backhed, F., Fulton, L. & Gordon, J. I. Diet-Induced Obesity Is Linked to Marked but Reversible Alterations in the Mouse Distal Gut Microbiome. *Cell Host Microbe* **3**, 213–223 (2008). DOI 10.1016/j.chom.2008.02.015. [arXiv](#).
11. Ridaura, V. K. *et al.* Gut microbiota from twins discordant for obesity modulate metabolism in mice. *Sci.* **341** (2013). DOI 10.1126/science.1241214.
12. Ramezani, A. & Raj, D. S. The Gut Microbiome, Kidney Disease, and Targeted Interventions. *J. Am. Soc. Nephrol.* **25**, 657–670 (2014). URL <http://www.jasn.org/cgi/doi/10.1681/ASN.2013080905>. DOI 10.1681/ASN.2013080905.
13. Michail, S. *et al.* Alterations in the gut microbiome of children with severe ulcerative colitis. *Inflamm. Bowel Dis.* **18**, 1799–1808 (2012). DOI 10.1002/ibd.22860.
14. Gevers, D. *et al.* The treatment-naïve microbiome in new-onset Crohn’s disease. *Cell Host Microbe* **15**, 382–392 (2014). DOI 10.1016/j.chom.2014.02.005.
15. Irwin, J. *et al.* A rolling phenotype in Crohn’s disease. *PLoS ONE* **12** (2017). DOI 10.1371/journal.pone.0174954.
16. Pascal, V. *et al.* A microbial signature for Crohn’s disease. *Gut* **66**, 813–822 (2017). DOI 10.1136/gutjnl-2016-313235. [NIHMS150003](#).
17. Kappelman, M. D. *et al.* The Prevalence and Geographic Distribution of Crohn’s Disease and Ulcerative Colitis in the United States. *Clin. Gastroenterol. Hepatol.* **5**, 1424–1429 (2007). DOI 10.1016/j.cgh.2007.07.012.
18. Gilbert, J. A. & Neufeld, J. D. Life in a World without Microbes. *PLoS Biol.* **12** (2014). DOI 10.1371/journal.pbio.1002020.
19. Pinto, A. J., Xi, C. & Raskin, L. Bacterial community structure in the drinking water microbiome is governed by filtration processes. *Environ. Sci. Technol.* **46**, 8851–8859 (2012). DOI 10.1021/es302042t.
20. Moran, M. A. The global ocean microbiome (2015). DOI 10.1126/science.aac8455.
21. Sunagawa, S. *et al.* Structure and function of the global ocean microbiome. *Sci.* **348** (2015). DOI 10.1126/science.1261359. [NIHMS150003](#).
22. Fierer, N. Embracing the unknown: Disentangling the complexities of the soil microbiome (2017). DOI 10.1038/nrmicro.2017.87.
23. Chaparro, J. M., Sheflin, A. M., Manter, D. K. & Vivanco, J. M. Manipulating the soil microbiome to increase soil health and plant fertility (2012). DOI 10.1007/s00374-012-0691-4.
24. Hamady, M. & Knight, R. Microbial community profiling for human microbiome projects: Tools, techniques, and challenges (2009). DOI 10.1101/gr.085464.108. [0402594v3](#).
25. Janda, J. M. & Abbott, S. L. 16S rRNA gene sequencing for bacterial identification in the diagnostic laboratory: Pluses, perils, and pitfalls (2007). DOI 10.1128/JCM.01228-07.
26. Chakravorty, S., Helb, D., Burday, M., Connell, N. & Alland, D. A detailed analysis of 16S ribosomal RNA gene segments for the diagnosis of pathogenic bacteria. *J. Microbiol. Methods* **69**, 330–339 (2007). DOI 10.1016/j.mimet.2007.02.005.
27. Caporaso, J. G. *et al.* QIIME allows analysis of high-throughput community sequencing data (2010). DOI 10.1038/nmeth.f.303. [NIHMS150003](#).
28. Lawley, B. & Tannock, G. W. *Analysis of 16S rRNA Gene Amplicon Sequences Using the QIIME Software Package*, vol. 1537 (Springer, 2017). URL <http://www.ncbi.nlm.nih.gov/pubmed/27924593>{%}0Ahttp://link.springer.com/10.1007/978-1-4939-6685-1{__}9.
29. Schloss, P. D. *et al.* Introducing mothur: Open-source, platform-independent, community-supported software for describing and comparing microbial communities. *Appl. Environ. Microbiol.* **75**, 7537–7541 (2009). DOI 10.1128/AEM.01541-09.
30. Edgar, R. C., Haas, B. J., Clemente, J. C., Quince, C. & Knight, R. UCHIME improves sensitivity and speed of chimera detection. *Bioinforma.* **27**, 2194–2200 (2011). DOI 10.1093/bioinformatics/btr381.
31. Knights, D., Costello, E. K. & Knight, R. Supervised classification of human microbiota (2011). DOI 10.1111/j.1574-6976.2010.00251.x.

32. Statnikov, A. *et al.* A comprehensive evaluation of multicategory classification methods for microbiomic data. *Microbiome* **1**, 11 (2013). URL <http://microbiomejournal.biomedcentral.com/articles/10.1186/2049-2618-1-11>. DOI 10.1186/2049-2618-1-11.
33. Xu, X. *et al.* Metadp: a comprehensive web server for disease prediction of 16s rRNA metagenomic datasets. *Biophys. Reports* **2**, 106–115 (2016).
34. Eck, A. *et al.* Robust microbiota-based diagnostics for inflammatory bowel disease. *J. Clin. Microbiol.* **55**, 1720–1732 (2017). DOI 10.1128/JCM.00162-17.
35. Duvallat, C., Gibbons, S. M., Gurry, T., Irizarry, R. A. & Alm, E. J. Meta-analysis of gut microbiome studies identifies disease-specific and shared responses. *Nat. Commun.* **8**, 1784 (2017). URL <http://www.nature.com/articles/s41467-017-01973-8>. DOI 10.1038/s41467-017-01973-8.
36. Carrieri, A. P. *et al.* *Host Phenotype Prediction from Differentially Abundant Microbes Using RoDEO*, 27–41 (Springer International Publishing, Cham, 2017). URL https://doi.org/10.1007/978-3-319-67834-4_3.
37. Cordier, T. *et al.* Predicting the Ecological Quality Status of Marine Environments from eDNA Metabarcoding Data Using Supervised Machine Learning. *Environ. Sci. & Technol.* **51**, 9118–9126 (2017). URL <http://pubs.acs.org/doi/abs/10.1021/acs.est.7b01518>. DOI 10.1021/acs.est.7b01518.
38. Fierer, N. *et al.* Forensic identification using skin bacterial communities. *Proc. Natl. Acad. Sci. United States Am.* **107**, 6477–81 (2010). URL <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2852011&tool=pmcentrez&rendertype=abstract%5Cnhttp://www.pnas.org/content/107/14/6477.short>. DOI 10.1073/pnas.1000162107.
39. Schmedes, S. E. *et al.* Targeted sequencing of clade-specific markers from skin microbiomes for forensic human identification. *Forensic Sci. Int. Genet.* **32**, 50–61 (2018).
40. Pasolli, E., Truong, D. T., Malik, F., Waldron, L. & Segata, N. Machine Learning Meta-analysis of Large Metagenomic Datasets: Tools and Biological Insights. *PLoS Comput. Biol.* **12** (2016). DOI 10.1371/journal.pcbi.1004977.
41. Costello, E. K. *et al.* Bacterial community variation in human body habitats across space and time. *Sci. (New York, N.Y.)* **326**, 1694–7 (2009). URL <http://www.sciencemag.org/content/326/5960/1694.short%5Cnhttp://dx.doi.org/10.1126/science.1177486%5Cnhttp://www.ncbi.nlm.nih.gov/pubmed/19892944%5Cnhttp://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC3602444>. DOI 10.1126/science.1177486.
42. Dutilh, B. E. *et al.* Explaining microbial phenotypes on a genomic scale: Gwas for microbes. *Briefings Funct. Genomics* **12**, 366–0380 (2013). DOI 10.1093/bfpg/elt008.
43. Ross, E. M., Moate, P. J., Marett, L. C., Cocks, B. G. & Hayes, B. J. Metagenomic Predictions: From Microbiome to Complex Health and Environmental Phenotypes in Humans and Cattle. *PLoS ONE* **8** (2013). DOI 10.1371/journal.pone.0073056.
44. Cui, H. & Zhang, X. Alignment-free supervised classification of metagenomes by recursive SVM. *BMC Genomics* **14** (2013). DOI 10.1186/1471-2164-14-641.
45. Asgari, E. & Mofrad, M. R. Continuous distributed representation of biological sequences for deep proteomics and genomics. *PLoS ONE* **10** (2015). DOI 10.1371/journal.pone.0141287. 1503.05140.
46. Min, S., Lee, B. & Yoon, S. Deep learning in bioinformatics. *Briefings Bioinforma.* bbw068 (2016). URL <https://academic.oup.com/bib/article-lookup/doi/10.1093/bib/bbw068>. DOI 10.1093/bib/bbw068. 1603.06430.
47. Ditzler, G., Polikar, R. & Rosen, G. Multi-Layer and Recursive Neural Networks for Metagenomic Classification. *IEEE Transactions on Nanobioscience* **14**, 608–616 (2015). DOI 10.1109/TNB.2015.2461219.
48. Nguyen, N.-P., Warnow, T., Pop, M. & White, B. A perspective on 16S rRNA operational taxonomic unit clustering using sequence similarity. *npj Biofilms Microbiomes* **2**, 16004 (2016). URL <http://www.nature.com/articles/npjbiofilms20164>. DOI 10.1038/npjbiofilms.2016.4.
49. McDonald, D. *et al.* An improved Greengenes taxonomy with explicit ranks for ecological and evolutionary analyses of bacteria and archaea. *ISME J.* **6**, 610–618 (2012). DOI 10.1038/ismej.2011.139.
50. Quast, C. *et al.* The SILVA ribosomal RNA gene database project: Improved data processing and web-based tools. *Nucleic Acids Res.* **41** (2013). DOI 10.1093/nar/gks1219.
51. Rideout, J. R. *et al.* Subsampled open-reference clustering creates consistent, comprehensive OTU definitions and scales to billions of sequences. *PeerJ* **2**, e545 (2014). URL <https://peerj.com/articles/545>. DOI 10.7717/peerj.545.
52. Cai, Y. *et al.* ESPRIT-Forest: Parallel clustering of massive amplicon sequence data in subquadratic time. *PLoS Comput. Biol.* **13** (2017). DOI 10.1371/journal.pcbi.1005518.
53. Koeppel, A. F. & Wu, M. Surprisingly extensive mixed phylogenetic and ecological signals among bacterial Operational Taxonomic Units. *Nucleic Acids Res.* **41**, 5175–5188 (2013). DOI 10.1093/nar/gkt241.

54. He, Y. *et al.* Erratum to: Stability of operational taxonomic units: an important but neglected property for analyzing microbial diversity. *Microbiome* **3**, 34 (2015). URL <http://www.microbiomejournal.com/content/3/1/34>. DOI 10.1186/s40168-015-0098-1.
55. Marçais, G. & Kingsford, C. A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinforma.* **27**, 764–770 (2011). DOI 10.1093/bioinformatics/btr011. [1006.1266v2](https://doi.org/10.1093/bioinformatics/btr011).
56. McHardy, A. C., Martín, H. G., Tsirigos, A., Hugenholtz, P. & Rigoutsos, I. Accurate phylogenetic classification of variable-length DNA fragments. *Nat. Methods* **4**, 63–72 (2007). DOI 10.1038/nmeth976.
57. Patil, K. R. *et al.* Taxonomic metagenome sequence assignment with structured output models (2011). DOI 10.1038/nmeth0311-191. [NIHMS150003](https://doi.org/10.1038/nmeth0311-191).
58. Wood, D. E. & Salzberg, S. L. Kraken: Ultrafast metagenomic sequence classification using exact alignments. *Genome Biol.* **15** (2014). DOI 10.1186/gb-2014-15-3-r46. [/www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3006164&tool=pmcentrez&rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3006164&tool=pmcentrez&rendertype=abstract).
59. Kawulok, J. & Deorowicz, S. CoMeta: Classification of metagenomes using k-mers. *PLoS ONE* **10** (2015). DOI 10.1371/journal.pone.0121453.
60. Menzel, P. & Krogh, A. Kaiju : Fast and sensitive taxonomic classification for metagenomics. *bioRxiv* **7**, 1–9 (2015). URL <http://dx.doi.org/10.1038/ncomms11257>. DOI 10.1101/031229.
61. Vervier, K., Mahé, P., Tournoud, M., Veyrieras, J. B. & Vert, J. P. Large-scale machine learning for metagenomics sequence classification. *Bioinforma.* **32**, 1023–1032 (2016). DOI 10.1093/bioinformatics/btv683. [1505.06915v1](https://doi.org/10.1093/bioinformatics/btv683).
62. Peterson, J. *et al.* The NIH Human Microbiome Project. *Genome Res.* **19**, 2317–2323 (2009). DOI 10.1101/gr.096651.109.
63. Huttenhower, C. & Human Microbiome Project Consortium. Structure, function and diversity of the healthy human microbiome. *Nat.* **486**, 207–14 (2012). URL <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3564958&tool=pmcentrez&rendertype=abstract>. DOI 10.1038/nature11234. [NIHMS150003](https://doi.org/10.1038/nature11234).
64. Yang, C. C. & Iwasaki, W. MetaMetaDB: A database and analytic system for investigating microbial habitability. *PLoS ONE* **9** (2014). DOI 10.1371/journal.pone.0087126.
65. Olson, R. S., Cava, W. L., Mustahsan, Z., Varik, A. & Moore, J. H. *Data-driven advice for applying machine learning to bioinformatics problems*, 192–203 (WORLD SCIENTIFIC, 2017). URL http://www.worldscientific.com/doi/abs/10.1142/9789813235533_0018.
66. Suykens, J. A. K. & Vandewalle, J. Least Squares Support Vector Machine Classifiers. *Neural Process. Lett.* **9**, 293–300 (1999). URL <http://dx.doi.org/10.1023/A:1018628609742>. DOI 10.1023/A:1018628609742. [1018628609742](https://doi.org/10.1023/A:1018628609742).
67. Breiman, L. Random Forests. *Mach. Learn.* **45**, 5–32 (2001). DOI 10.1023/A:1010933404324. [/dx.doi.org/10.1023/A:1010933404324](http://dx.doi.org/10.1023/A:1010933404324).
68. Goodfellow, I., Bengio, Y. & Courville, A. *Deep Learning* (MIT Press, 2016). <http://www.deeplearningbook.org>.
69. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. Dropout: prevent NN from overfitting. *J. Mach. Learn. Res.* **15**, 1929–1958 (2014). URL <http://jmlr.org/papers/v15/srivastava14a.html> <http://jmlr.org/papers/volume15/srivastava14a/srivastava14a.pdf>. DOI 10.1214/12-AOS1000. [1102.4807](https://doi.org/10.1214/12-AOS1000).
70. Kingma, D. P. & Ba, J. L. Adam: a Method for Stochastic Optimization. *Int. Conf. on Learn. Represent. 2015* 1–15 (2015). DOI <http://doi.acm.org.ezproxy.lib.ucf.edu/10.1145/1830483.1830503>. [1412.6980](https://doi.org/10.1145/1830483.1830503).
71. Pedregosa, F. & Varoquaux, G. *Scikit-learn: Machine learning in Python*, vol. 12 (ACM, 2011). URL <http://dl.acm.org/citation.cfm?id=2078195>. [arXiv:1201.0490v2](https://arxiv.org/abs/1201.0490v2).
72. Jolliffe, I. T. & Jolliffe, I. T. *Principal Component Analysis* (1986). URL <http://www.google.com/search?client=safari&rlz=en-us&hl=en&oeq=Principal+Component+Analysis&oeq=UTF-8&oeq=UTF-8>.
73. Van Der Maaten, L. J. P. & Hinton, G. E. Visualizing high-dimensional data using t-sne. *J. Mach. Learn. Res.* **9**, 2579–2605 (2008). URL http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=pubmed&cmd=Retrieve&dopt=AbstractPlus&list_uids=7911431479148734548related:VOiAgwMNY20J. DOI 10.1007/s10479-011-0841-3. [1307.1662](https://doi.org/10.1007/s10479-011-0841-3).
74. Bengio, Y., Courville, A. & Vincent, P. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis Mach. Intell.* **35**, 1798–1828 (2013). DOI 10.1109/TPAMI.2013.50. [1206.5538](https://doi.org/10.1109/TPAMI.2013.50).
75. Richard, F. J., Poulsen, M., Drijfhout, F., Jones, G. & Boomsma, J. J. Specificity in chemical profiles of workers, brood and mutualistic fungi in *Atta*, *Acromyrmex*, and *Sericomyrmex* fungus-growing ants. *J. Chem. Ecol.* **33**, 2281–2292 (2007). DOI 10.1007/s10886-007-9385-z.