



Universiteit
Leiden
The Netherlands

Metagenomics : beyond the horizon of current implementations and methods

Khachatryan, L.

Citation

Khachatryan, L. (2020, April 28). *Metagenomics : beyond the horizon of current implementations and methods*. Retrieved from <https://hdl.handle.net/1887/87513>

Version: Publisher's Version

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/87513>

Note: To cite this publication please use the final published version (if applicable).

Cover Page



Universiteit Leiden



The handle <http://hdl.handle.net/1887/87513> holds various files of this Leiden University dissertation.

Author: Khachatryan, L.

Title: Metagenomics : beyond the horizon of current implementations and methods

Issue Date: 2020-04-28

Reference-free resolving of long-read metagenomic data

L. Khachatryan¹, S. Y. Anvar¹, R. H. A. M. Vossen³, and J. F. J. Laros^{1,2,4}

1 Department of Human Genetics, Leiden University Medical Center, Leiden, The Netherlands

2 Clinical Genetics, Leiden University Medical Center, Leiden, The Netherlands

3 Leiden Genome Technology Center, Leiden University Medical Center, Leiden, The Netherlands

4 GenomeScan, Leiden, The Netherlands

bioRxiv 2019 <https://doi.org/10.1101/811760>

3.1 Background

The analysis of metagenomic data is becoming a routine for many different research fields, since it serves scientific purposes as well as improves our life quality. Particularly, with the use of metagenomics a large step was made towards the understanding of the human microbiome and uncovering its real composition and diversity [225, 226, 227, 228, 229, 230]. The understanding of the human microbiome in health and disease contributed to the development of diagnostics and treatment strategies based on metagenomics knowledge [231, 232, 233, 234, 235, 236, 237, 238]. The study of microbial ecosystems allows us to predict the possible processes, changes and sustainability of particular environments [239, 240]. Genes isolated from uncultivable inhabitants of soil metagenomes are being successfully utilized, for example, in the biofuel industry for production and tolerance to byproducts [30, 241, 242]. Various newly discovered biosynthetic capacities of microbial communities benefit the manufacturing of industrial, food, and health products, as well as contribute into the field of bioremediation [54, 55, 56, 57].

Despite all the progress made in resolving genetic data derived from environmental samples, it is still a challenging task. Reads binning is one of the most critical steps in the analysis of metagenomic data. To estimate the composition of a particular microbiome, it is important to ensure that sequencing reads derived from the same organism are grouped together. Currently, alignment of DNA extracted from an environmental sample to a set of known sequences remains the main strategy for metagenomics binning [243, 244]. There is a full range of techniques allowing the comparison of metagenomic reads to a reference database. It can be performed using different metagenomic data types (16S or WGS) and various matching approaches (classic alignment or matching performed using k -mers or taxonomic signatures). Most of the time, the binning is performed for all reads in the database, but in some cases only a particular subset of sequencing data is selected for binning. Lastly, there is a wide spectrum of databases that can be used to perform the binning. The database might contain all possible annotated nucleotide/protein sequences, marker genes for distinct phylogenetic clades, sequencing signatures specific to particular taxa, etc. (see more detailed explanation in Chapter 1). The obvious downside of all listed strategies is the incapability to perform an accurate binning for the reads derived from organisms that are not present in the reference database.

Metagenomics binning was improved by alignment-free approaches, which can be split into two subgroups: reference-dependent and reference-independent methods. The tools from the first subgroup utilize existing databases to train a supervised classifier for the reads binning. Various techniques can be performed to achieve this goal: Support Vector Machine, Interpolated Markov Models, Gaussian Mixture Models, Hidden Markov Models [147, 148, 149, 151, 152, 153, 150]. Even though these approaches are reference dependent, they can be used to classify reads derived

from previously unknown species. However, the accuracy of reference-dependent methods will be always limited by the content of reference databases. The content of the current reference databases utilized for training differs from the true distribution of microbial species on our planet [245, 246, 247, 248, 249, 250, 251]. For some metagenomic datasets the amount of unknown sequences might be quite high [252, 253], thus using supervised classification tools based on known genetic sequences is questionable if this is the case.

Reference-independent approaches for metagenomics binning try to solve the problem of missing taxonomic content: they are designed to classify reads into genetically homogeneous groups without utilizing any information from known genomes. Instead, they use only the features of the sequencing data (usually k -mer distributions, DNA segments of length k) for classification. One of those tools, LickelyBin, performs a Markov Chain Monte Carlo approach based on the assumption that the k -mer frequency distribution is homogeneous within a bacterial genome [140]. This tool performs well for very simple metagenomes with significant phylogenetic diversity within the metagenome, but it cannot handle genomes with more complicated structure such as those resulting from horizontal gene transfer [141]. Another one, AbundanceBin [142], works under the assumption that the abundances of species in metagenome are following a Poisson distribution, and thus struggles analysing datasets where some species have similar abundance ratios. MetaCluster [143] and BiMeta [144] address this problem of non-Poisson species distribution. However, for these tools it is necessary to provide an estimation of the final number of clusters, which cannot be done for many metagenomes without any prior knowledge. Also, both MetaCluster and BiMeta are using a Euclidian metric to compute the dissimilarity between k -mer profiles, which was shown to be influenced by stochastic noise in analysed sequences [145]. Another recent tool, MetaProb, implements a more advanced similarity measure technique and can automatically estimate the number of read clusters [146]. This tool classifies metagenomic datasets in two steps: first, reads are grouped based on the extent of their overlap. After that, a set of representing reads is chosen for each group. Based on the comparison of the k -mer distributions for those sets, groups are merged together into final clusters. Even though MetaProb outperformed other tools during the analysis of simulated data, it was shown to perform not very well on the real metagenomic data.

In this article we present a new technique for alignment- and reference-free classification of metagenomic data. Our approach is based on a pairwise comparison of k -mer profiles calculated for each sequencing read in a long-read metagenomic dataset, using the previously described kPAL toolkit [213]. It also performs unsupervised clustering to facilitate the identification of genetically homogeneous groups of reads present in a sample. The main assumption of our method is that after assigning the pairwise distances for all reads in the dataset, those belonging to the same organism will form dense groups, and thus the metagenome binning could be resolved using

density-based clustering. We developed an algorithm which automatically detects the regions with high density and hierarchically splits the dataset until there is one dense region per cluster. The approach is designed to work with long reads (more than 1,000 bp) since we calculate k -mer profiles for each read separately and shorter reads would yield non-informative profiles. We performed our analysis on long PacBio reads that were either simulated or generated from a real metagenomic sample. We have shown that despite the fact that PacBio data is known to have a high error rate, the approach successfully performed read classification for simulated and real metagenomic data.

3.2 Materials and Methods

3.2.1 Software

All analyses were done using publicly available tools (parameters used are listed below for each specific case) along with custom Python scripts which are stored in a Git repository¹.

3.2.2 PacBio data simulation

Complete genomes of five common skin bacteria were used to generate artificial PacBio metagenomes (see Table 3.1). The reads were simulated from reference sequences using the PBSIM toolkit [254] with CLR as the output data type and a final sequencing depth of 20. For the calibration of the read length distribution, a set of previously sequenced *C. difficile* reads [255] was used as a model.

3.2.3 Bioreactor metagenome PacBio sequencing

Bioreactor metagenome coupling anaerobic ammonium oxidation (Annamox) to Nitrite/Nitrate dependent Anaerobic Methane Oxidation (N-DAMO) processes [256] was used to generate WGS PacBio sequencing data.

Metagenome contained the N-DAMO bacteria *Methylomirabilis oxyfera* (complete genome with GeneBank Accession FP565575.1 was used as a reference), two Annamox bacteria (*Kuenenia stuttgartiensis*, assembly contigs from the Bio Project PRJEB22746 were used as a reference and a member of *Broccardia* genus, assembly contigs of *Broccardia sinica* from Bio Project PRJDB103 were used as reference) and an archaea species *Methanoperedens nitroreducens* (assembly contigs from the Bio Project PRJNA242803 were used as a reference).

Bacterial cell pellets were disrupted with a Dounce homogenizer. DNA was isolated using a Genomic Tip 500/G kit (Qiagen) and needle sheared with a 26G blunt end

¹Available at <https://git.lumc.nl/1.khachatryan/pacbio-meta>

needle (SAI Infusion). Pulsed-field Gel electrophoresis was performed to assess the size distribution of the sheared DNA. A SMRTbell library was constructed using 5 μ g of DNA following the 20 kb template preparation protocol (Pacific Biosciences). The SMRTbell library was size selected using the BluePippin system (SAGE Science) with a 10 kb lower cut-off setting. The final library was sequenced with the P6-C4 chemistry with a movie time of 360 minutes.

3.2.4 Reads origin checking

Reads were corrected using the PacBio Hierarchical Genome Assembly Process algorithm before being mapped to the genomes of the expected metagenome inhabitants genomes using the BLASR aligner [257] with default settings. The alignments were used to determine the origin of the reads. Reads that were not mapped during the previous step were subjected to the BLASTn [102] search against the NCBI database. The identity cut-off was set to 90, the (E)value was chosen to be 0.001.

3.2.5 Bioreactor metagenome PacBio reads assembly

The assembly of corrected PacBio reads was performed using the FALCON [258] assembler. The resulting contigs were mapped to the candidate reference genomes using LAST [104] with default settings. To determine the similarity cutoff for the mapping procedure, the curve representing the number of contigs versus the similarity to the reference genome was analysed. The first inflection point at (in case of mapping contigs to the *M. oxyfera* genome 12%), dividing the fast-declining part of the curve from the slow-declining part, was chosen as a threshold (See Supplementary materials for more details).

3.2.6 Binning procedure

For each read, the frequencies of all possible five-mers were calculated using the *count* command of the kPAL toolkit. The resulting profiles were balanced (a procedure that compensates for differences that occur because of reading either the forward or reverse complement strand) and compared in a pairwise manner by using the *balance* and *matrix* commands of kPAL accordingly, yielding a pairwise distance matrix. Normalization for differences in read length was dealt with by the scaling option during the pairwise comparison.

The resulting distance matrix, hereafter called the original distance matrix, was subjected to a multi-step clustering procedure. A schematic representation of this procedure can be found in Figure 3.1. Due to practical limitations (runtime), this analysis was restricted to a set of 10,000 randomly selected reads. This multi-step clustering procedure works recursively: it starts with the analysis of a set of reads and either reports the entire set as one cluster, or it splits the set into two subsets,

which are each analysed using the same procedure. The decision whether to split the set of reads into two subsets is made using the following approach. First, the pairwise distances for all reads in the set are extracted from the original distance matrix in order to construct the working distance matrix. After that, the dimensionality of the analysed set is decreased to three using the t-SNE algorithm [259] in order to reduce noise caused by outliers in the distance matrix. The reads, now represented by a point in three-dimensional space, are subjected to density-based clustering using the DBSCAN algorithm [260] with the default distance function.

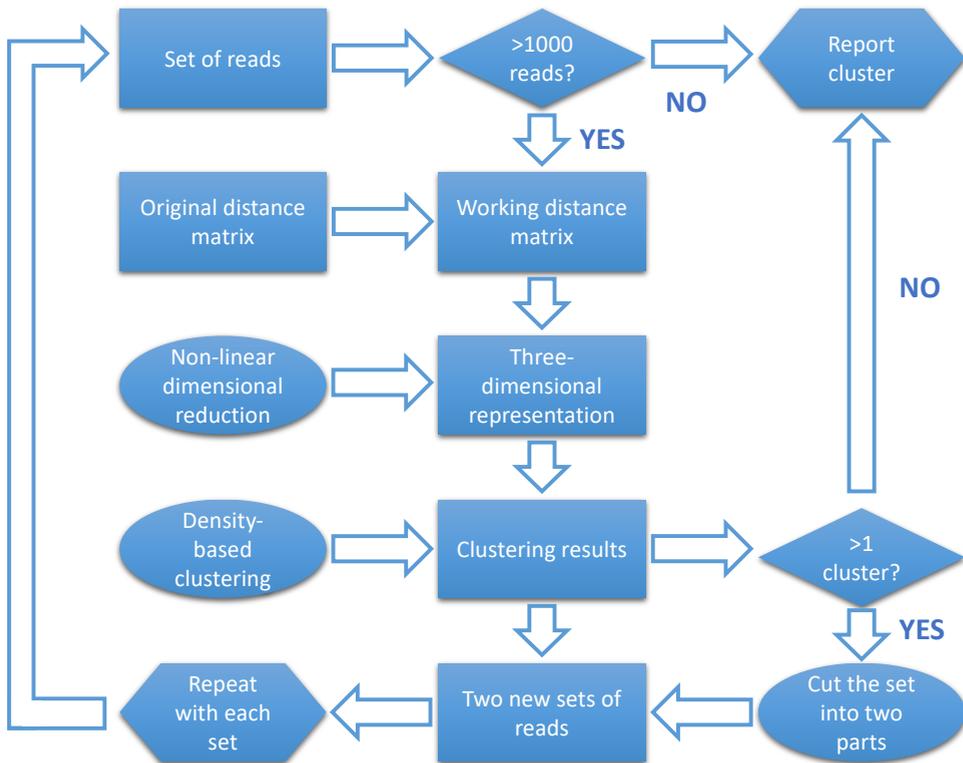


Figure 3.1: Schematic representation of the clustering procedure.

We choose the *MinPts* parameter of DBSCAN (the minimal amounts of points in the neighborhood to extend the cluster) to be either 1% of the size of the dataset for sets larger than 2,000 reads, or 20 for sets smaller than 2,000 reads. The number of clusters found by DBSCAN depends on the neighborhood diameter ϵ . When ϵ is too small, no clusters are reported since all points are isolated. On the other hand, when ϵ is too large all points are grouped into one cluster. Our algorithm

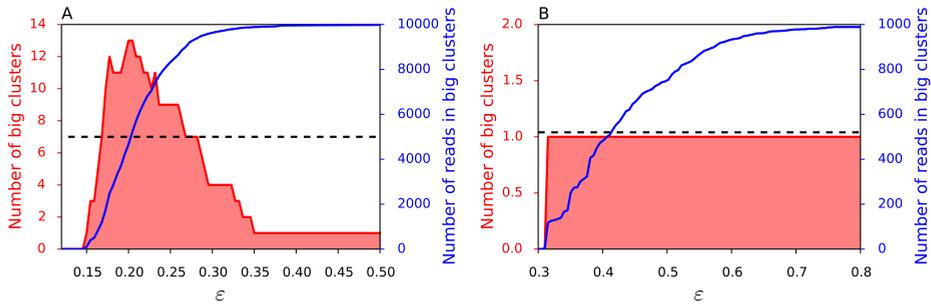


Figure 3.2: Density-based clustering analysis example. The data is clustered with DBSCAN with ϵ ranging from 0 to the value when 90% of the points are assigned to one cluster. When at least half of the data set is assigned to a dense cluster, the number of clusters is used to determine whether subdivision of the data set is required. Only if more than one cluster is identified at this point, the procedure is repeated recursively with two partitions of the data. The partitions are determined by using the largest ϵ that clusters the data into two clusters. In this example two datasets are shown: one that was further split into two partitions (A) and one that was reported as one dense cluster (B).

therefore performs a parameter sweep for ϵ , from the value providing zero clusters to the value with which 99% of the reads are grouped in one cluster for the chosen *MinPts*. The results of this parameter sweep are used to check the dependency of the number of dense clusters on a particular ϵ (only clusters larger than 100 points are considered) and how many points of the analysed set are included in the obtained clusters (Figure 3.2). If for some ϵ there are two or more clusters that together cover more than half of the total amount, the analysed set is divided into two new sets (Figure 3.2A). The analysed set is reported as one cluster if the aforementioned condition is not satisfied (Figure 3.2B), or when the size of the analysed set is smaller than 1,000 points.

The division is done using the following strategy. DBSCAN is performed using the optimal ϵ , yielding two dense clusters that serve as center points for two partitions. Each of the remaining unclassified points is assigned to the cluster containing the closest classified neighbor.

3.2.7 Classification for larger sets

Read classification for sets larger than 10,000 was performed in two steps. First, 10,000 reads (larger than 10 kb) were randomly chosen and classified using the algorithm described in section 3.2.6. After that, the pairwise distances between every unclassified read and every classified read were calculated using their 5-mer profiles. These distances were used to assign the unclassified read to the cluster containing the closest classified read.

3.2.8 Data availability

Sequencing reads of bioreactor metagenome were submitted to SRA under the study number SRP159147.

Supplementary materials were deposited on Figshare and available for downloading using the following link: <https://doi.org/10.6084/m9.figshare.c.4218857.v1>

3.3 Results

3.3.1 Reads classification in artificial PacBio metagenomes

To construct artificial metagenomes, we used simulated PacBio reads based on the genomes of five common skin flora bacteria together with so-called "noise" reads. These are reads from a PacBio sequencing data of an environmental metagenome [261] that were not assigned to the major inhabitant *K. stuttgartiensis* or other known organisms. They were added to represent low abundant species that are present in any typical metagenomic dataset.

We constructed four artificial PacBio datasets in this way, each containing 10,000 randomly selected reads (length > 9 kb) containing 0%, 5%, 10% and 15% noise reads, respectively. For the simplicity the number of simulated reads was adjusted to provide an equal abundance for each bacterium in the final metagenome (Table 3.1).

Reads origin	RefSeq AC	Genome length, Mb	Number of reads per dataset			
			0%	5%	10%	15%
<i>S. mitis</i>	NC_013853.1	2.1	1,246	1,183	1,121	1,059
<i>P. acnes</i>	NC_017550.1	2.5	1,443	1,371	1,298	1,226
<i>S. epidermidis</i>	NC_004461.1	2.6	1,448	1,376	1,304	1,231
<i>A. calcoaceticus</i>	NC_016603.1	3.9	2,236	2,125	2,013	1,901
<i>P. aeruginosa</i>	NC_002516.2	6.3	3,627	3,446	3,264	3,083

Table 3.1: Content of artificial metagenomics PacBio datasets.

We subjected each dataset to the classification procedure described in section 3.2.6. The reads in the resulting clusters were then classified according to their origin (See Supplementary Material for more data). In Figure 3.3, it can be seen that for each experiment we obtained five large clusters (> 1,000 reads) consisting mainly of reads belonging to the same species. For all three datasets containing noise reads we see the tendency of noise reads to be clustered with some fraction of *P. acnes* and *P. aeruginosa* reads. However, as can be seen from Figure 3.3 and Table 3.2, increasing the noise content leads to better isolation of these reads. Indeed, for dataset B (5% of the noise reads), the majority of noise reads were assigned to the cluster that is primarily occupied by reads belonging to *P. acnes* and *P. aeruginosa*. Increasing the noise content (dataset C and D in Fig. 4, 10% and 15% noise reads accordingly) led to the appearance of two clusters which contain mostly noise reads (Table 3.2, A). We also see that with the increase of noise content, the fractions of *P. acnes* and *P. aeruginosa* reads included in the same clusters as the noise reads are dropping (Table 3.2, B). In conclusion, the more noise reads were added to the dataset, the better they were grouped together in one or two clusters (Table 3.2, A).

Dataset Reads origin	5% noise	10% noise		15% noise	
	Cluster 2	Cluster 2	Cluster 8	Cluster 6	Cluster 7
A					
noise	21.4	90.3	47.8	85.6	97.3
<i>P. acnes</i>	63.7	0.5	33.8	5.6	0
<i>P. aeruginosa</i>	10.4	1.3	19.1	8.9	0
B					
noise	91.8	55.9	39.9	45.0	50.8
<i>P. acnes</i>	99.6	0.2	22.3	3.6	0
<i>P. aeruginosa</i>	6.4	0.2	5.3	2.3	0

Table 3.2: Composition of clusters containing the majority of noise reads after the classification procedure for three artificial PacBio datasets. A - cluster composition; B - the percentage of reads with particular origin (noise, *P. acnes* or *P. aeruginosa*) included to the cluster within all reads of the same origin in the dataset. Clusters are grouped per dataset. Only organisms whose reads would occupy more than 90% of cluster content are shown.

3.3.2 PacBio sequencing of bioreactor metagenome

After sequencing and correction, we obtained 31,757 reads longer than 1kb for the bioreactor metagenome. The read length distribution for this dataset can be found in Figure 3.4. Reads were mapped to the genomes of the expected metagenome inhabitants. Since the groups of reads that we could map to the genomes of *K. stuttgartiensis* and *B. sinica* had a significant overlap (27%), we decided to combine reads mapped to the reference genomes of these two organisms in one group. We detected almost no (0.01%) reads that would map to the *M. nitroreducens* genome in the sequencing data, suggesting that this organism was either not present in the metagenome sample, or that its DNA could not be isolated reliably during the sample preparation.

Thus, we divided our reads into three groups: uniquely mapped on *M. oxyfera* (4,903 reads), uniquely mapped on *K. stuttgartiensis*/*B. sinica* (2,973 reads), and all remaining reads with unknown origin (75%, 23,881 reads). The reads with unknown origin were checked with the BLASTn software against NCBI microbial database, to find significant similarity to any known organism. However, only 334 reads (less than 2% of total number of checked reads) got hits; there were no organisms among the obtained hits reported more than 53 times.

3.3.3 Bioreactor metagenome PacBio read classification

For the reads originating from *M. oxyfera* and *K. stuttgartiensis*/*B. sinica*, we checked whether the data was clustered by origin. Since roughly 75% of this sequencing data is of unknown origin, we assessed whether the clustering results for reads with unknown origin is robust. To do this, we created five subsets using the bioreactor

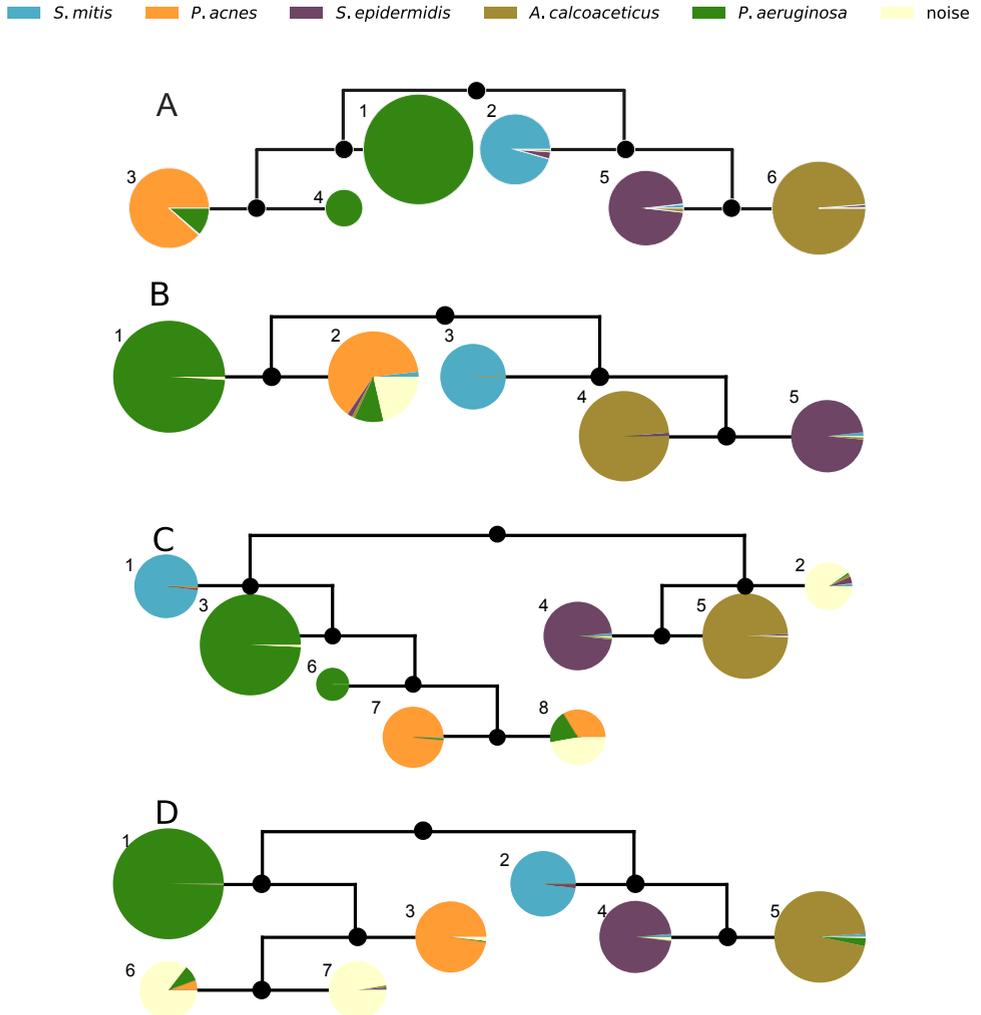


Figure 3.3: Classification recall for artificial PacBio metagenomes. Subsets that were subjected to the partitioning are shown as black circles, final clusters are represented as pie charts with the colour indicating the reads origin. The area of the pie chart corresponds to the relative cluster size. The cluster number is shown next to each pie chart. The results are shown for datasets with 0% (A), 5% (B), 10% (C) and 15% (D) of noise reads.

metagenome sequencing data. Each subset contains 10,000 randomly selected reads with length > 10 kb. After subjecting each subset to the classification procedure, we checked whether reads, shared by two subsets, are being clustered similarly. We compared all clusters from different subsets in a pairwise manner and marked two

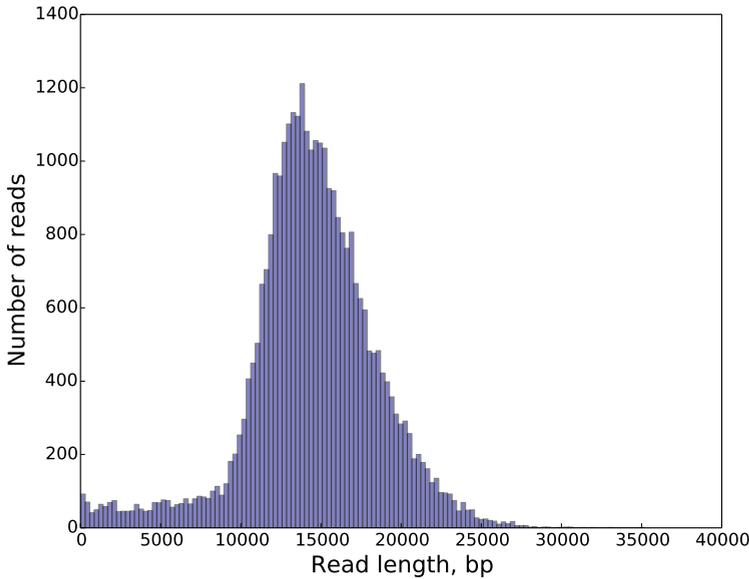


Figure 3.4: Bioreactor metagenome reads length distribution.

clusters 'similar' when they shared at least 25% of their content. On average, every pair of subsets shared 34% of their content. Thus, in case of perfect matching of clustering results, the pair of clusters from two different subsets should on average share 34% of their content. The 25% cutoff value was chosen to compensate for possible flaws introduced by clustering mis-assignments.

In Figure 3.5 this analysis is shown as a graph: each pie chart represents a cluster obtained for one of the subsets (with a subset number marked next to the pie chart). Clusters are connected if they were marked as similar and thus shared more than 25% of their content. We looked for sub-graphs, of size five for which all five nodes would be mutually connected. That would mean that all five clusters are coming from the different subsets and share a significant (at least 25% out of 34% possible) number of reads. These groups of clusters (here and after called the stable groups) represent reads that are clustered the same way regardless of the subset of reads selected. Clusters belonging to the stable groups are called the stable clusters. The proportion of reads in the stable clusters was comparable among datasets and equaled on average 64%. As displayed in Figure 3.5, we found seven groups of stable clusters. Four groups of stable clusters have clusters with more than 1,000 reads, and two of those four are represented by clusters enriched with *M. oxyfera* or *K. stuttgartiensis*/*B. sinica* reads. In Table 3.3 we display the content and the number of reported clusters after the classification procedure for each of the five subsets.

Once we estimated the robustness of the classification procedure, we selected the

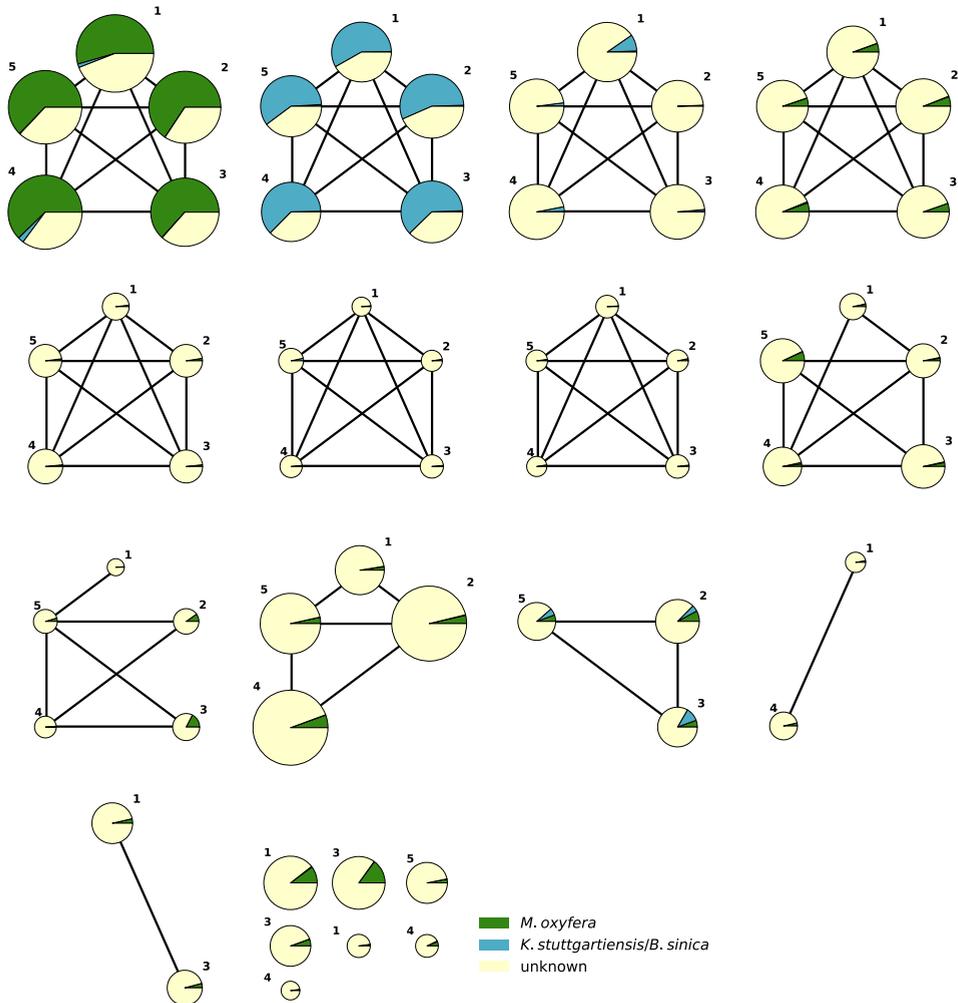


Figure 3.5: Comparison of classification results obtained for five bioreactor sub-datasets. The pie charts represent reported clusters for all sub-datasets coloured by the origin of reads in cluster. The pie chart area indicates the relative size of the cluster. The number next to the node denotes the sub-dataset, for which the cluster was obtained. Two clusters are connected with a node if they belong to two different sub-datasets and share at least 25% of their content. The groups of size five (the set of five fully connected pie-charts) represent groups of stable clusters.

Subset	1	2	3	4	5
Number of <i>M. oxyfera</i> reads	1,499	1,563	1,528	1,544	1,529
Number of <i>K. stuttgartiensis</i> / <i>B. sinica</i> reads	949	918	981	935	906
Clusters after the classification procedure	14	11	13	13	12
Big (>1,000 reads) clusters	5	5	5	5	5
% of reads in stable clusters	65.96	64.12	61.98	64.46	64.16

Table 3.3: Subsets information and clustering results.

subset that yielded the lowest number of clusters (subset 2, 11 clusters) for downstream analysis. The content of all clusters that were not reported as stable were merged into one cluster. Thus, the original 10,000 reads were spread among 8 clusters. These clusters were used as a classifier for the remaining 21,757 reads in the dataset (Table 3.4).

Cluster	Stable	Reads before extension	Reads after extension
1	Yes	403	1,038
2	Yes	168	528
3	Yes	1,133	3,204
4	Yes	1,540	5,151
5	Yes	1,004	3,337
6	Yes	181	506
7	Yes	1,983	6,459
8	No	3,588	11,534

Table 3.4: Results of bioreactor metagenome reads classification

3.3.4 Assembly of the bioreactor metagenome before and after reads binning

We assembled reads belonging to different clusters separately, and compared the resulting contigs with the results of the assembly of the entire dataset. The total number of contigs after assembly of the partitioned dataset was comparable to the amount of contigs obtained from the assembly of the entire dataset (Table 3.5). The same can be said about the total length of contigs and contigs length distributions (see supplementary materials). These results, showing that the database partitioning did not lead to the change of the contigs number or their lengths, can be seen as indirect evidence proving that our *k*-mer based binning of metagenome reads results in species-based clustering.

We compared the assembled contigs obtained for the entire and partitioned datasets to the reference genomes of *M. oxyfera*, *K. stuttgartiensis* and *B. sinica*. Even though

we could successfully map around 9% of the reads to the reference genomes of *K. stuttgartiensis* and *B. sinica*, we did not get contigs that could be mapped to these genomes. However, the contigs assembled from the entire and partitioned datasets did map to *M. oxyfera* genome. Only 91 out of 196 contigs obtained from the entire dataset assembly could be mapped back to the *M. oxyfera* genome covering 54% of its length. For the assembly of the partitioned dataset, 85 contigs were mapped to the genome of *M. oxyfera* in total, covering 52.65% of its length. The vast majority of those contigs (79, covering 51% of the *M. oxyfera* genome length) derived from the assembly of reads belonging to one cluster. Thus, our dataset partitioning binned the majority of contigs according to their origin.

Dataset assembled	Entire dataset	Cl 1	Cl 2	Cl 3	Cl 4	Cl 5	Cl 6	Cl 7	Cl 8
Assembly length, bp	3,251,357	5,438	10,747	380,905	377,792	601,065	0	1,602,878	41,310
Contigs	196	1	1	28	30	47	0	71	2
Contigs mapped on <i>M. oxyfera</i> genome	91	0	0	9	1	2	0	71	2
Length of mapped contigs	1,842,182	0	0	132,863	11,945	21,105	0	1,497,132	17,013
<i>M. oxyfera</i> genome covered, %	54	0	0	1.2	0.1	0.15	0	51	0

Table 3.5: Results of entire and partitioned bioreactor sequencing data assembly and comparison of obtained contigs to the *M. oxyfera* genome. Cl - cluster.

3.4 Discussion

We described a new approach for efficient, alignment-free binning of metagenomic sequencing reads based on k -mer frequencies. Our method successfully classifies reads per organism of origin, for both simulated and real metagenomic data.

As shown in the results section, the approach was used to classify reads obtained by

PacBio sequencing of a real bioreactor metagenome. The absolute majority of the reads with known origin (*M. oxyfera* or *K. stuttgartiensis*/*B. sinica*) were clustered together per origin after pairwise comparison of their k mer profiles and subsequent density-based cluster detection. This result was robust, as we observed during the analysis of five subsets of the original PacBio sequencing data with overlapping content. The same experiment demonstrated that each subset provides a similar number of clusters. Reads with unknown origin had a tendency to cluster similarly among different subsets, again confirming the clustering consistency. Although the majority of reads in the analysed metagenome was of unknown origin, the results can be used to estimate the microbial community complexity for its most abundant inhabitants.

The binning of the bio-reactor metagenomic dataset had almost no influence on the results of the metagenome assembly. The number of contigs and their lengths obtained for the entire and partitioned datasets were comparable. This indicates that the k -mer based reads binning leads to the organism-based partitioning of metagenomic data. Furthermore, contigs, belonging to the same organism, were automatically grouped together when assembling the dataset subjected to the classification procedure. Thus, our k -mer based binning technique can be used to interpret metagenomic assembly results.

Performing the binning procedure on an artificially generated PacBio datasets lead to a reads classification per organism, even after adding reads with unknown origin (noise reads). Moreover, increasing the proportion of noise reads leads to a better separation between them and the reads with known origin. This observation supports the ck -mer central hypothesis of this research, namely that k -mer distances can be used to cluster reads of the same origin together once those reads provide sufficient coverage of the organisms' genome.

The main disadvantages of the current implementation of our method is the limited number of reads (10,000) that can be analysed. As mentioned before, reads, derived from the same organism, will cluster together, but this is possible only under the condition that the organisms' genome is sufficiently covered. Thus, the described technique is unsuitable for the analysis of metagenomes with a large number of inhabitants or when the inhabitants have large genomes, as 10,000 reads will not be enough to provide sufficient coverage. The depth of the classification that can be performed by the suggested method is still to be discovered.

We believe that adapting our metagenomics reads binning technique for larger sets of data and further investigation of its metagenome resolving capacity would allow to expand the current limits of microbiology in the future.

3.5 Author Statements

3.5.1 Funding information

This research is financed by a grant number 727.011.002 of the Netherlands Organisation for Scientific Research (NWO).

3.5.2 Acknowledgements

We would like to thank the group of Prof. Huub Op den Camp for the bioreactor metagenome material, Prof. Boudewijn P. F. Lelieveldt for the idea to perform dimensional reduction using t-SNE, and Martijn Vermaat for the help with coding.

3.5.3 Conflicts of interest

The authors declare that there are no conflicts of interest.

