# Reference-free resolution of long-read metagenomic data

Lusine Khachatryan[1*], Seyed Yahya Anvar[1,2†], Rolf H. A. M. Vossen[2††], and Jeroen F. J. Laros[1,3,4†††]

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

2 - Leiden Genome Technology Center, Leiden University Medical Center, Leiden, The Netherlands

3 - Department of Clinical Genetics, Leiden University Medical Center, Leiden, The Netherlands

4 - GenomeScan, Leiden, The Netherlands


* - khachatryan.l.a@gmail.com – Corresponding author

† - s.y.anvar@lumc.nl

†† - r.h.a.m.vossen@lumc.nl

††† - j.f.j.laros@lumc.nl

18    ABSTRACT

19    *Background*

20    Read binning is a key step in proper and accurate analysis of metagenomics data. Typically, this is

21    performed by comparing metagenomics reads to known microbial sequences. However, microbial

22    communities usually contain mixtures of hundreds to thousands of unknown bacteria. This restricts

23    the accuracy and completeness of alignment-based approaches. The possibility of reference-free

24    deconvolution of environmental sequencing data could benefit the field of metagenomics,

25    contributing to the estimation of metagenome complexity, improving the metagenome assembly, and

26    enabling the investigation of new bacterial species that are not visible using standard laboratory or

27    alignment-based bioinformatics techniques.

28    *Results*

29    Here, we apply an alignment-free method that leverages on $k$-mer frequencies to classify reads within

30    a single long read metagenomic dataset. In addition to a series of simulated metagenomic datasets,

31    we generated sequencing data from a bioreactor microbiome using the PacBio RSII single-molecule

32    real-time sequencing platform. We show that distances obtained after the comparison of $k$-mer

33    profiles can reveal relationships between reads within a single metagenome, leading to a clustering

34    per species.

35    *Conclusions*

36    In this study, we demonstrated the possibility to detect substructures within a single metagenome

37    operating only with the information derived from the sequencing reads. The obtained results are

38    highly important as they establish a principle that might potentially expand the toolkit for the detection

39    and investigation of previously unknow microorganisms.

40

41    KEYWORDS

42    Metagenomics binning, PacBio sequencing, metagenome resolving

43

44  INTRODUCTION

45  The analysis of metagenomics data is becoming a routine for many different research fields, since it

46  serves scientific purposes as well as improves our life quality. Particularly, with the use of

47  metagenomics a large step was made towards the understanding of the human microbiome and

48  uncovering its real composition and diversity [1-6]. The understanding of the human microbiome in

49  health and disease contributed to the development of diagnostics and treatment strategies based on

50  metagenomic knowledge [7-14]. The study of microbial ecosystems allows us to predict the possible

51  processes, changes and sustainability of particular environments [15, 16]. Genes isolated from

52  uncultivable inhabitants of soil metagenomes are being successfully utilized, for example, in the

53  biofuel industry for production and tolerance to byproducts [17-19]. Various newly discovered

54  biosynthetic capacities of microbial communities benefit the production of industrial, food, and health

55  products, as well as contribute into the field of bioremediation [20-23].

56  Despite all the progress made in resolving genetic data derived from environmental samples, it is still

57  a challenging task. Reads binning is one of the most critical steps in the analysis of metagenomics

58  data. To estimate the composition of a particular microbiome, it is important to ensure that sequencing

59  reads derived from the same organism are grouped together. Currently, alignment of DNA extracted

60  from an environmental sample to a set of known sequences remains the main strategy for

61  metagenomics binning [24, 25]. There is a full range of techniques allowing the comparison of

62  metagenomic reads to a reference database. It can be performed using different metagenomic data

63  types (16S or WGS) and various matching approaches (classic alignment or use of $k$-mers or

64  taxonomical signatures). Most of the time, the binning is performed for all reads in the database, but

65  in some cases only a particular subset of sequencing data is selected for binning. Lastly, there is a

66  wide spectrum of databases that can be used to perform the binning. The database might contain all

67  possible annotated nucleotide/protein sequences, marker genes for distinct phylogenetic clades,

68  sequencing signatures specific to particular taxa, etc. The obvious downside of all listed strategies is

3

69    the incapability to perform an accurate binning for the reads of organisms that are not present in the

70    reference database.

71    Metagenomic binning was improved by alignment-free approaches, which can be split into two

72    subgroups: reference-dependent and reference-independent methods. The tools from the first

73    subgroup utilize existing databases to train a supervised classifier for the reads binning. Various

74    techniques can be performed to achieve this goal: linear regression, Interpolated Markov Models,

75    Gaussian Mixture Models, Hidden Markov Models [26-32]. Even though these approaches are

76    reference dependent, they can be used to classify reads that are derived from previously unknown

77    species. However, the accuracy of reference-dependent methods will be always limited by the content

78    of reference databases. The content of the current reference databases utilized for training differs from

79    the true distribution of microbial species on our planet [33-39]. For some metagenomic datasets the

80    amount of unknown sequences might be quite high [40, 41], thus using supervised classification tools

81    based on known genetic sequences is questionable in such cases.

82    Reference-independent approaches for metagenomics binning try to solve the problem of missing

83    taxonomic content: they are designed to classify reads into genetically homogeneous groups without

84    utilizing any information from known genomes. Instead, they use only the features of the sequencing

85    data (usually $k$-mer distributions, DNA segments of length $k$) for classification. One of those tools,

86    LicklyBin, performs a Markov Chain Monte Carlo approach based on the assumption that the $k$-mer

87    frequency distribution is homogeneous within a bacterial genome [42]. This tool performs well for

88    very simple metagenomes with significant phylogenetic diversity within the metagenome, but it

89    cannot handle genomes with more complicated structure such as those resulting from horizontal gene

90    transfer [43]. Another one, AbundanceBin [44], works under the assumption that the abundances of

91    species in metagenome reads are following a Poisson distribution, and thus struggles analyzing

92    datasets where some species have similar abundance ratios. MetaCluster [45] and BiMeta [46] address

93    this problem of non-Poisson species distribution. However, for these tools it is necessary to provide

94    an estimation of the final number of clusters, which cannot be done for many metagenomes without

95    any prior knowledge. Also, both MetaCluster and BiMeta are using a Euclidian metric to compute the

96    dissimilarity between $k$-mer profiles, which was shown to be influenced by stochastic noise in

97    analyzed sequences [47]. Another recent tool, MetaProb, implements a more advanced similarity

98    measure technique and can automatically estimate the number of read clusters [48]. This tool classifies

99    metagenomic datasets in two steps: first, reads are grouped based on the extent of their overlap. After

100    that, a set of representing reads is chosen for each group. Based on the comparison of the $k$-mer

101    distributions for those sets, groups are merged together into final clusters. Even though MetaProb

102    outperformed other tools during the analysis of simulated data, it was shown to perform not very well

103    on the real metagenomics data.

104    In this article we present a new technique for alignment- and reference-free classification of

105    metagenomics data. Our approach is based on a pairwise comparison of $k$-mer profiles calculated for

106    each sequencing read in a long-read metagenomics dataset, using the previously described kPAL

107    toolkit [49]. It also performs unsupervised clustering to facilitate the identification of genetically

108    homogeneous groups of reads present in a sample. The main assumption of our method is that after

109    assigning the pairwise distances for all reads in the dataset, those belonging to the same organism will

110    form dense groups, and thus the metagenome binning could be resolved using density-based

111    clustering. We developed an algorithm which automatically detects the regions with high density and

112    hierarchically splits the dataset until there is one dense region per cluster. The approach is designed

113    to work with long reads (more than 1000 bp) since we calculate $k$-mer profiles for each read separately

114    and shorter reads would yield non-informative profiles. We performed our analysis on long PacBio

115    reads that were either simulated or generated from a real metagenomic sample. We have shown that

116    despite the fact that PacBio data is known to have a high error rate, the approach successfully

117    performed read classification for simulated and real metagenomic data.

118    MATERIALS AND METHODS

119    *1. Software*

120    All analyses were done using publicly available tools (parameters used are listed below for each

121    specific case) along with custom Python scripts.

122    *2. PacBio data simulation*

123    Complete genomes of five common skin bacteria were used to generate artificial PacBio

124    metagenomes (Table 1). The reads were simulated from reference sequences using the PBSIM toolkit

125    [50] with CLR as the output data type and a final sequencing depth of 20. For the calibration of the

126    read length distribution, a set of previously sequenced *C. difficille* reads [51] was used as a model.

127    *3. Bioreactor metagenome PacBio sequencing*

128    Bioreactor metagenome coupling anaerobic ammonium oxidation (Annamox) to Nitrite/Nitrate

129    dependent Anaerobic Methane Oxidation (N-DAMO) processes [52] was used to generate WGS

130    PacBio sequencing data.

131    Metagenome contained the N-DAMO bacteria *Methylomirabilis oxyfera* (complete genome with

132    GeneBank Acsession FP565575.1 was used as a reference), two Annamox bacteria (*Kuenenia*

133    *stuttgartiensis*, assembly contigs from the Bio Project PR- JEB22746 were used as a reference and a

134    member of *Broccardia* genus, assembly contigs of *Broccardia sinica* from Bio Project PRJDB103

135    were used as reference) and an archaea species *Methanoperedens nitroreducens* (assembly contigs

136    from the Bio Project PRJNA242803 were used as a reference).

137    Bacterial cell pellets were disrupted with a Dounce homogenizer. DNA was isolated using a Genomic

138    Tip 500/G kit (Qiagen) and needle sheared with a 26G blunt end needle (SAI Infusion). Pulsed-field

139    Gel electrophoresis was performed to assess the size distribution of the sheared DNA. A SMRTbell

140    library was constructed using $5\mu g$ of DNA following the 20kb template preparation protocol (Pacific

141    Biosciences). The SMRTbell library was size selected using the BluePippin system (SAGE Science)

142 with a 10kb lower cut-off setting. The final library was sequenced with the P6-C4 chemistry with a

143 movie time of 360 minutes.

144 *4. Reads origin checking*

145 Reads were corrected using the PacBio Hierarchical Genome Assembly Process algorithm before

146 being mapped to the genomes of the references of expected metagenome inhabitants using the BLASR

147 aligner [53] with default settings. The alignments were used to determine the origin of the reads.

148 Reads that were not mapped during the previous step were subjected to the BLASTn [54] search

149 against the NCBI database. The identity cut-off was set to 90, the (E)value was chosen to be 0.001.

150 *5. Bioreactor metagenome PacBio reads assembly*

151 The assembly of corrected PacBio reads was performed using the FALCON [55] assembler. The

152 resulting contigs were mapped to the candidate reference genomes using LAST [56] with default

153 settings. To determine the similarity cutoff for the mapping procedure, the curve representing the

154 number of contigs versus the similarity to the reference genome was analyzed. The first inflection

155 point at (in case of mapping contigs to the *M.oxyfera* genome 12%), dividing the fast-declining part

156 of the curve from the slow-declining part, was chosen as a threshold (See Section S1 of Additional

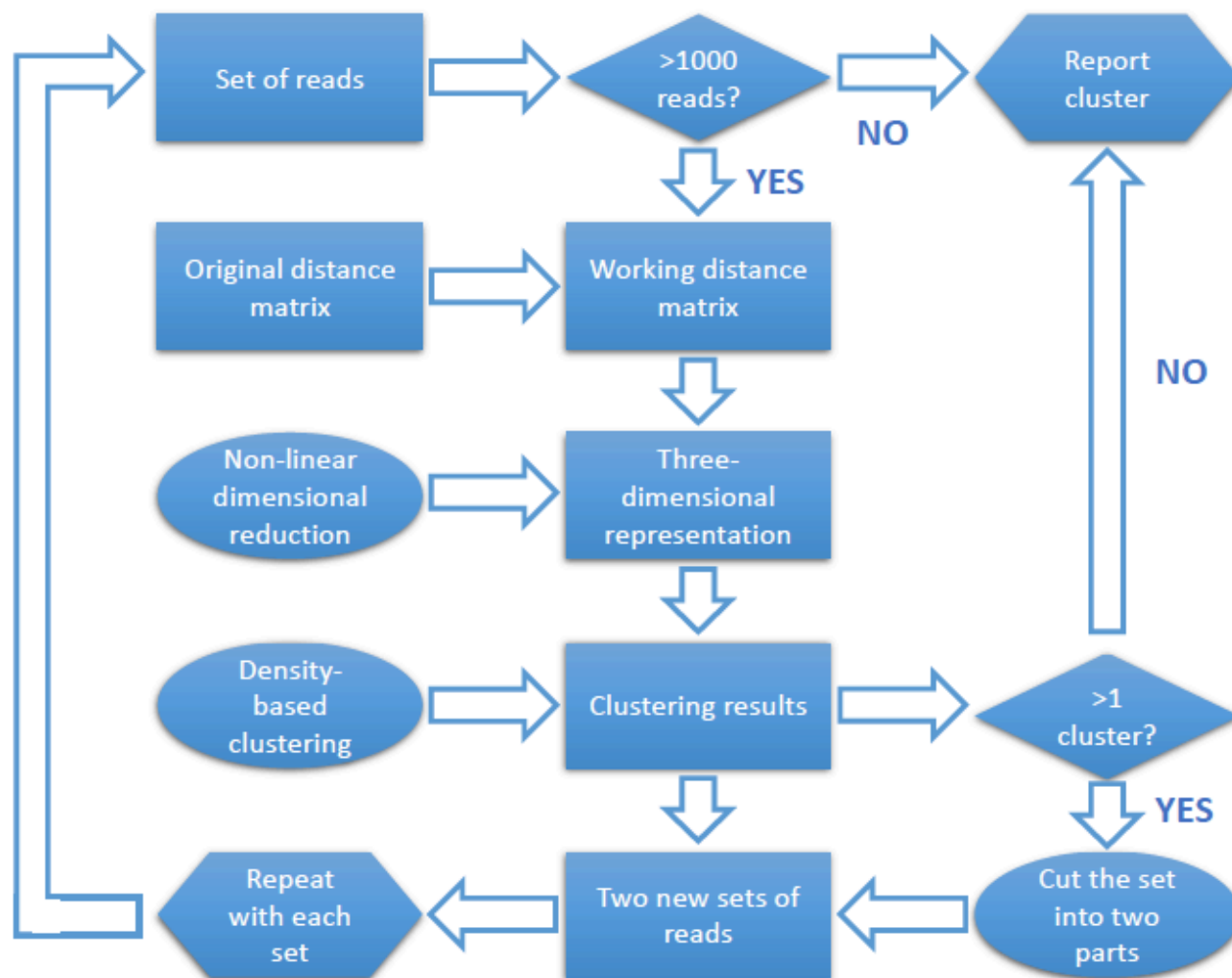157 file 1 for more details).

158 *6. Binning procedure*

159 For each read, the frequencies of all possible five-mers are calculated using the *count* command of

160 the kPAL toolkit. The resulting profiles are balanced (a procedure that compensates for differences

161 that occur because of reading either the forward or reverse complement strand) and compared in a

162 pairwise manner by using the *balance* and *matrix* commands of kPAL accordingly, yielding a pairwise

163 distance matrix. Normalization for differences in read length is dealt with by the scaling option during

164 the pairwise comparison.

165 The resulting distance matrix, hereafter called the original distance matrix, was subjected to a multi-

166 step clustering procedure. A schematic representation of this procedure can be found in Fig. 1. Due

7

167    to practical limitations (runtime), this analysis was restricted to a set of 10 000 randomly selected

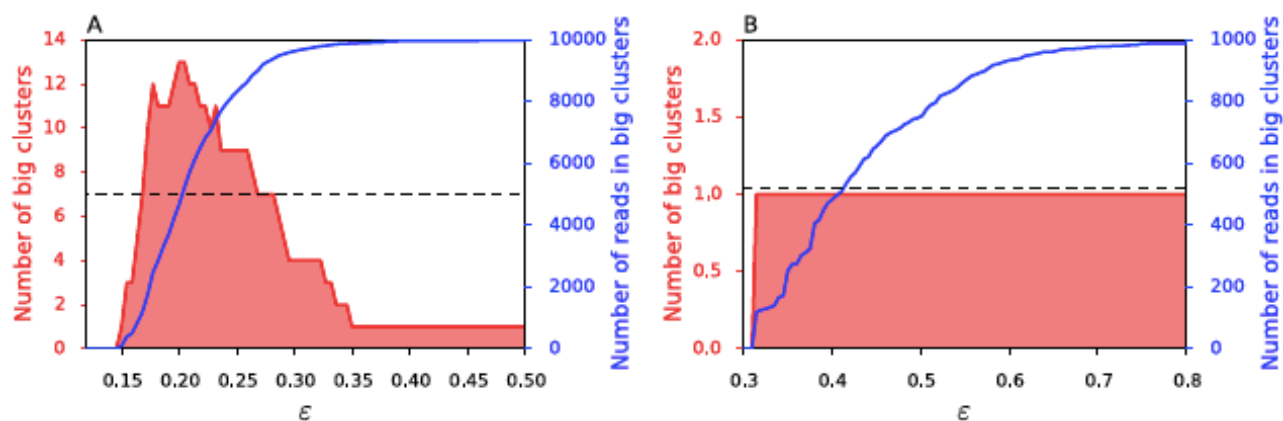168    reads.

169



170

171    Fig. 1. Schematic representation of the clustering procedure.

172

173    This multi-step clustering procedure works recursively: it starts with the analysis of a set of reads and

174    either reports the entire set as one cluster, or it splits the set into two subsets, which are each analyzed

175    using the same procedure. The decision whether to split the set of reads into two subsets is made using

176    the following approach. First, the pairwise distances for all reads in the set are extracted from the

177    original distance matrix in order to construct the working distance matrix. After that, the

8

178 dimensionality of the analyzed set is decreased to three using the t-SNE algorithm [57] in order to

179 reduce noise caused by outliers in the distance matrix. The reads, now represented by a point in three-

180 dimensional space, are subjected to density-based clustering using the DBSCAN algorithm [58] with

181 the default distance function. We choose the *MinPts* parameter of DBSCAN (the minimal amounts of

182 points in the neighborhood to extend the cluster) to be either 1% of the size of the dataset for sets

183 larger than 2000 reads, or 20 for sets smaller than 2000 reads. The number of clusters found by

184 DBSCAN depends on the neighborhood diameter $\varepsilon$. When $\varepsilon$ is too small, no clusters are reported

185 since all points are isolated. On the other hand, when $\varepsilon$ is too large all points are grouped into one

186 cluster. Our algorithm therefore performs a parameter sweep for $\varepsilon$, from the value providing zero

187 clusters to the value with which 99% of the reads are grouped in one cluster for the chosen *MinPts*.

188



190 Fig. 2. Density-based clustering analysis example. The data is clustered with DBSCAN with $\varepsilon$ ranging from 0 to the value

191 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

192 number of clusters is used to determine whether subdivision of the data set is required. Only if more than one cluster is

193 identified at this point, the procedure is repeated recursively with two partitions of the data. The partitions are determined

194 by using the largest $\varepsilon$ that clusters the data into two clusters. In this example two datasets are shown: one that was further

195 split into two partitions (A) and one that was reported as one dense cluster (B).

196

197 The results of this parameter sweep are used to check the dependency of the number of dense clusters

198 on a particular $\varepsilon$ (only clusters larger than 100 points are considered) and how many points of the

9

199     analyzed set are included in the obtained clusters (Fig. 2). If for some $\varepsilon$ there are two or more clusters

200     that together cover more than half of the total amount, the analyzed set is divided into two new sets

201     (Fig. 2A). The analyzed set is reported as one cluster if the aforementioned condition is not satisfied

202     (Fig. 2B), or when the size of the analyzed set was smaller than 1000 points.

203     The division is done using the following strategy. DBSCAN is performed using the optimal $\varepsilon$,

204     yielding two dense clusters that serve as center points for two partitions. Each of the remaining

205     unclassified points is assigned to the cluster containing the closest classified neighbor.

206     *7. Classification for larger sets*

207     Read classification for sets larger than 10 000 was performed in two steps. First, 10 000 reads (larger

208     than 10kb) were randomly chosen and classified using the algorithm described in previous section.

209     After that, the pairwise distances between every unclassified read and every classified read were

210     calculated using their 5-mer profiles. These distances were used to assign the unclassified read to the

211     cluster containing the closest classified read.

212     *8. Data availability*

213     Sequencing reads of bioreactor metagenome were submitted to NCBI under the BioProject number

214     PRJNA487927. Artificial PacBio metagenomic reads with the addition of $0\%, 5\%, 10\%,$ and $15\%$ of

215     real "noise" reads were submitted to NCBI under the BioProject number PRJNA533970.

216     Supplementary materials were deposited on Figshare and available for downloading using the

217     following link: https://doi.org/10.6084/m9.figshare.c.4218857.v1.

218     Example of the classification procedure can be found using the following link:

219     https://git.lumc.nl/l.khachatryan/pacbio-meta/blob/master/analysis/real_data/tsne_subset2/analysis_example.ipynb

220

10

221    RESULTS

222    *1. Reads classification in artificial PacBio metagenomes*

223    To construct artificial metagenomes, we used simulated PacBio reads based on the genomes of five

224    common skin flora bacteria together with so-called "noise" reads. These are reads from a PacBio

225    sequencing data of an environmental metagenome [59] that were not assigned to the major inhabitant

226    *K. stuttgartiensis* or other known organisms. They were added to represent low abundant species that

227    are present in any typical metagenomic dataset.

228    We constructed four artificial PacBio datasets in this way, each containing 10 000 randomly selected

229    reads (length > 9kb) containing 0%, 5%, 10% and 15% noise reads, respectively. For the simplicity

230    the number of simulated reads was adjusted to provide an equal abundance for each bacterium in the

231    final metagenome (see Table 1).

232    We subjected each dataset to the classification procedure described in Section 6 of MATERIALS

233    AND METHODS. The reads in the resulting clusters were then classified according to their origin

234    (See Section S2 of Additional file 1 for more data).

235

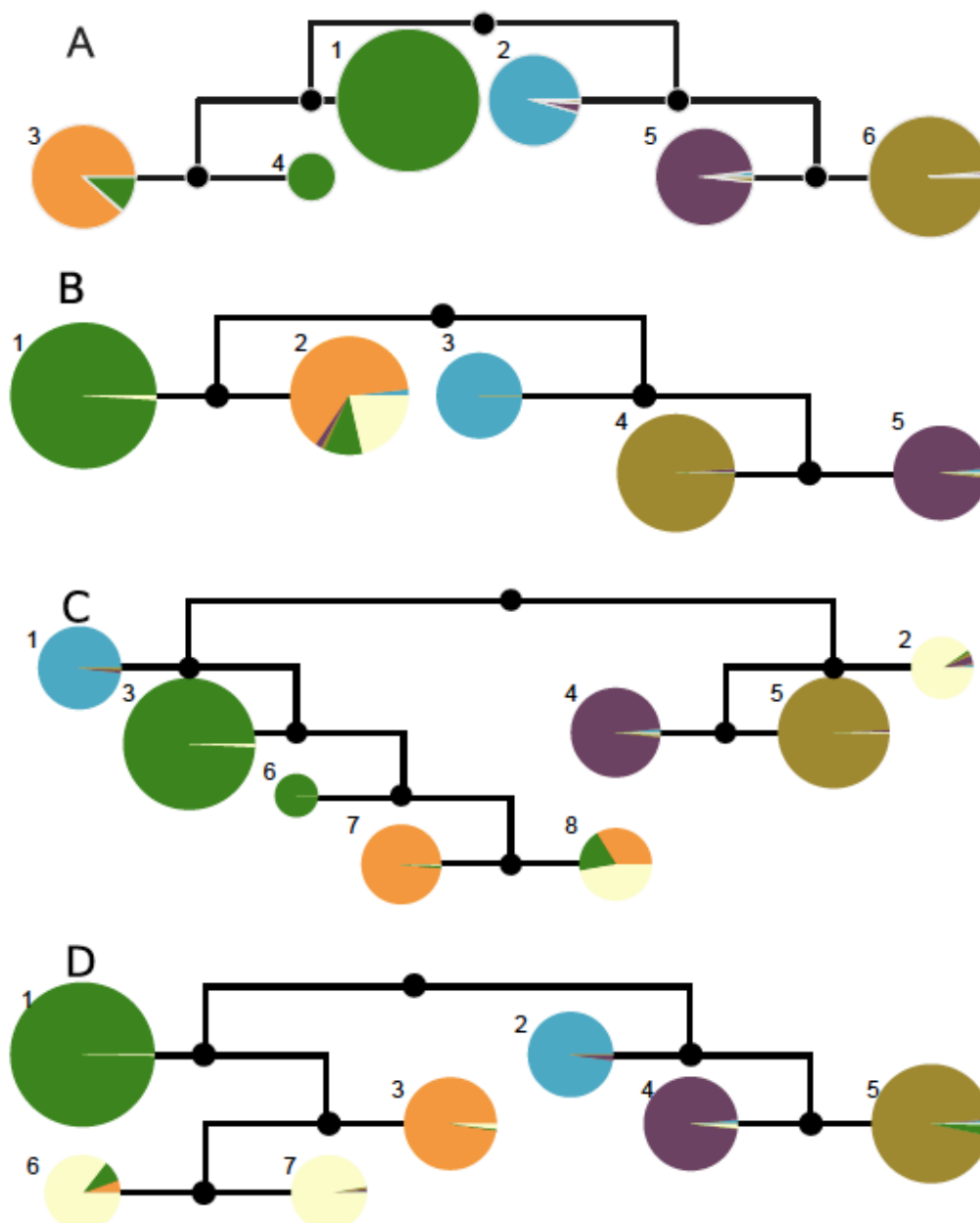236    **Table 1.** Content of artificial metagenomics PacBio datasets.

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

237

238    In Fig. 3, it can be seen that for each experiment we obtained five large clusters (> 1 000 reads)

239    consisting mainly of reads belonging to the same species.

240



241

242    Fig. 3. Classification recall for artificial PacBio metagenomes. Subsets that were subjected to the partitioning are shown

243    as black circles, final clusters are represented as pie charts with the color indicating the reads origin. The area of the pie

244    chart corresponds to the relative cluster size. The cluster number is shown next to each pie chart. The results are shown

245    for datasets with 0% (A), 5% (B), 10% (C) and 15% (D) of noise reads.

12

246    For all three datasets containing noise reads we see the tendency of noise reads to be clustered with

247    some fraction of *P. acnes* and *P. aeruginosa* reads.

248    However, as can be seen from Fig. 3 and Table 2, increasing the noise content leads to better isolation

249    of these reads. Indeed, for dataset B (5% of the noise reads), the majority of noise reads were assigned

250    to the cluster that is primarily occupied by reads belonging to *P. acnes and P. aeruginosa*. Increasing

251    the noise content (dataset C and D in Fig. 4, 10% and 15% noise reads accordingly) led to the

252    appearance of two clusters which contain mostly noise reads (Table 2, A).

253

254    **Table 2.** Composition of clusters containing the majority of noise reads after the classification procedure for three artificial

255    PacBio datasets.

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

256    A - cluster composition; B - the percentage of reads with particular origin (noise, *P. acnes* or *P. aeruginosa*) included to

257    the cluster within all reads of the same origin in the dataset. Clusters are grouped per dataset. Only organisms whose reads

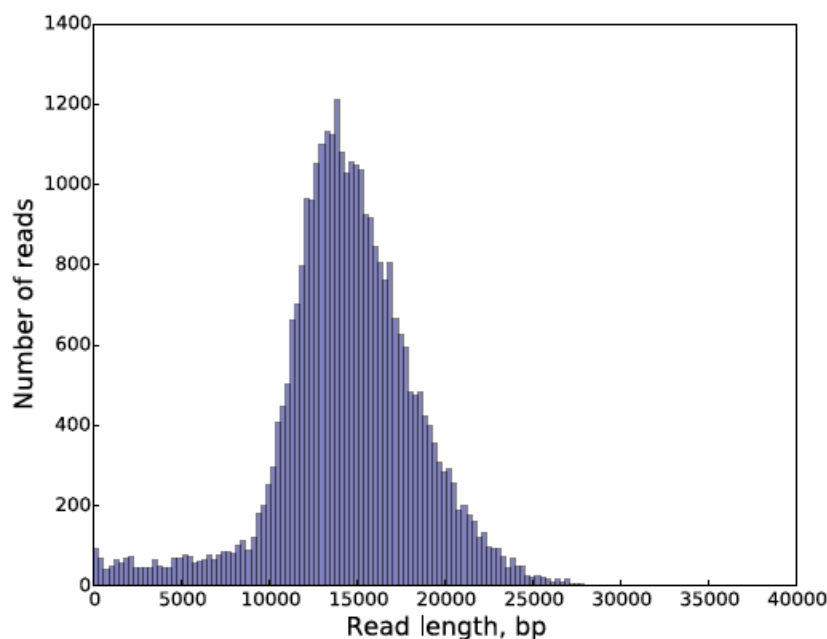258    would occupy more than 90% of cluster content are shown.

259

13

260    We also see that with the *increase* of noise content, the fractions of *P. acnes* and *P. aeruginosa* reads

261    included in the same clusters as the noise reads are dropping (Table 2, B). In conclusion, the more

262    noise reads were added to the dataset, the more they were grouped together in one or two clusters

263    (Table 2, A).

264    *4.2 PacBio sequencing of bio reactor metagenome*

265    After sequencing and correction, we obtained 31,757 reads longer than 1kb for the bio reactor

266    metagenome. The read length distribution for this dataset can be found in Fig. 4.

267



268

269    Fig. 4. Bio reactor metagenome reads length distribution.
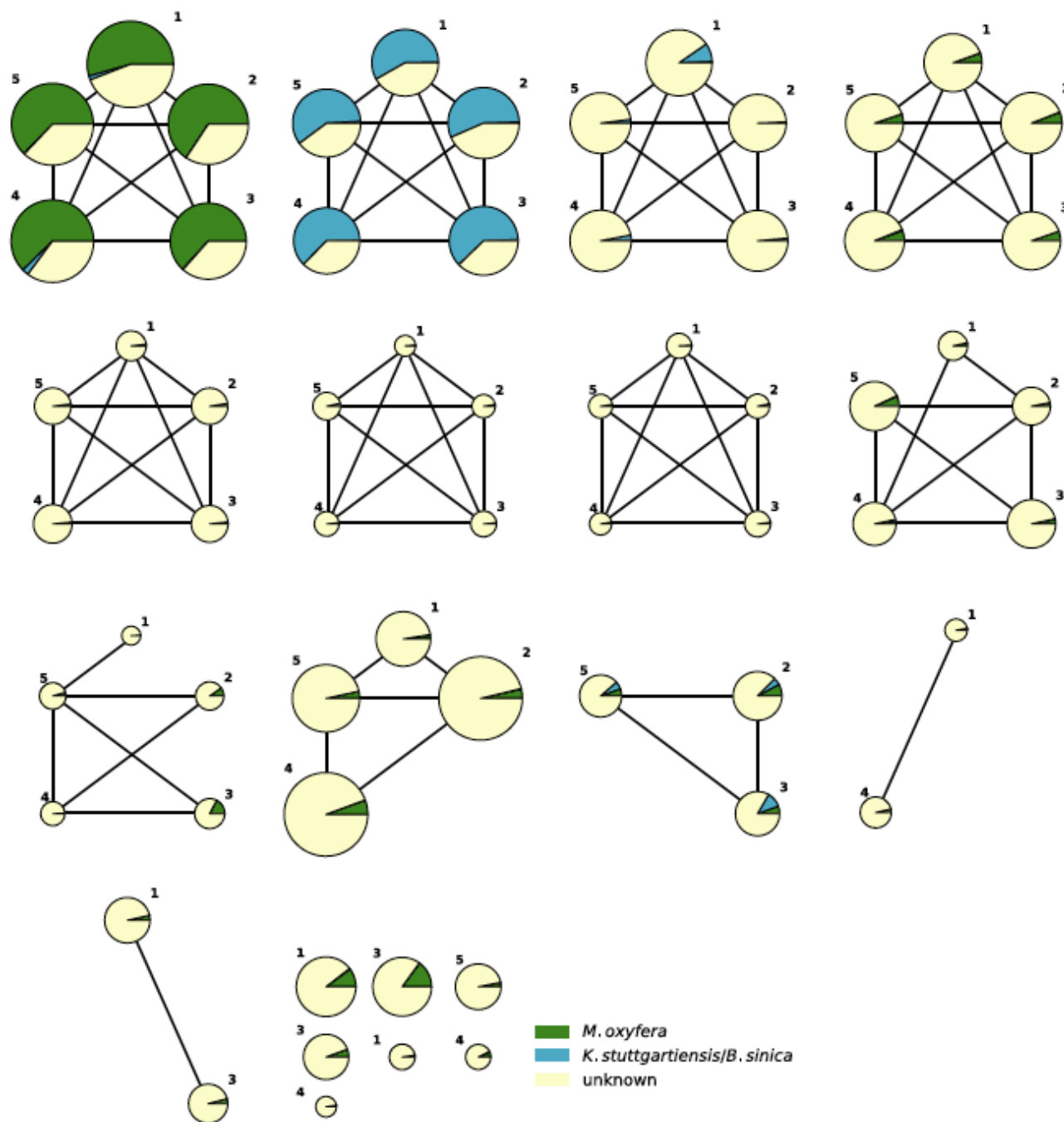
270

271    Reads were mapped to the genomes of the expected metagenome inhabitants or genomes of closely

272    related species. Since the groups of reads that we could map to the genomes of *K. stuttgartiensis* and

273    *B. sinica* had a significant overlap (27%), we decided to combine reads mapped to the reference

274    genomes of these two organisms in one group. We detected almost no (0.01%) reads that would map

275    to the *M. nitroreducens* genome in the sequencing data, suggesting that this organism was either not

276    present in the metagenome sample, or that its DNA could not be isolated reliably during the sample

14

277   preparation. Thus, we divided our reads into three groups: uniquely mapped on *M. oxyfera* (4,903

278   reads), uniquely mapped on *K. stuttgartiensis/B. sinica* (2973 reads), and all remaining reads with

279   unknown origin (~75%, 23881 reads). The reads with unknown origin were checked with the

280   BLASTn software against NCBI microbial database, to find significant similarity to any known

281   organism. However, only 334 reads (less then 2% of total number of checked reads) got hits; there

282   were no organisms among the obtained hits reported more than 53 times.

283   *4.3 Bio reactor metagenome PacBio read classification*

284   For the reads originating from *M. oxyfera* and *K. stuttgartiensis/B. sinica*, we checked whether the

285   data was clustered by origin. Since roughly 75% of this sequencing data is of unknown origin, we

286   assessed whether the clustering results for reads with unknown origin is robust. To do this, we created

287   five subsets using the bio reactor metagenome sequencing data. Each subset contains 10,000 randomly

288   selected reads with length > 10kb. After subjecting each subset to the classification procedure, we

289   checked whether reads, shared by two subsets, are being clustered similarly. We compared all clusters

290   from different subsets in a pairwise manner and marked two clusters 'similar' when they shared at

291   least 25% of their content. On average, every pair of subsets shared 34% of their content. Thus, in

292   case of perfect matching of clustering results, the pair of clusters from two different subsets should

293   on average share 34% of their content. The 25% cutoff value was chosen to compensate for possible

294   flaws introduced by clustering mis-assignments.  In Fig. 5 this analysis is shown as a graph: each pie

295   chart represents a cluster obtained for one of the subsets (with a subset number marked next to the pie

296   chart).

297

Fig. 5 Comparison of classification results obtained for five Bio reactor sub-datasets. The pie charts represent reported clusters for all sub-datasets colored 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.

305    Clusters are connected if they were marked as similar and thus shared more then 25% of their content.

306    We looked for sub-graphs, of size five for which all five nodes would be mutually connected. That

307    would mean that all five clusters are coming from the different subsets and share a significant (at least

308    25% out of 34% possible) number of reads. These groups of clusters (here and after called the stable

309    groups) represent reads that are clustered the same way regardless of the subset of reads selected.

310    Clusters belonging to the stable groups are called the stable clusters. The proportion of reads in the

311    stable clusters was comparable among datasets and equaled on average 64%. As displayed in Fig. 5,

312    we found seven groups of stable clusters. Four groups of stable clusters have clusters with more than

313    1 000 reads, and two of those four are represented by clusters enriched with *M. oxyfera* or

314    *K. stuttgartiensis/B. sinica* reads. In Table 3 we display the content and the number of reported

315    clusters after the classification procedure for each of the five subsets.

316

317    Table 3. Subsets information and clustering results.

| Subset | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| number of *M.oxyfera* reads | 1 499 | 1 563 | 1 528 | 1 544 | 1 529 |
| number of *K.stuttgartiensis/B.sinica* reads | 949 | 918 | 981 | 935 | 906 |
| Clusters after the classification procedure | 14 | 11 | 13 | 13 | 12 |
| Big (>1000 reads) clusters | 5 | 5 | 5 | 5 | 5 |
| % of reads in stable clusters | 65.96 | 64.12 | 61.98 | 64.46 | 64.16 |

318

319

320    Once we estimated the robustness of the classification procedure, we selected the subset that yielded

321    the lowest number of clusters (subset 2, 11 clusters) for downstream analysis. The content of all

322    clusters that were not reported as stable were merged into one cluster. Thus, the original 10 000 reads

17

323   were spread among 8 clusters. These clusters were used as a classifier for the remaining 21 757 reads

324   in the dataset (Table 4).

325

326   **Table 4.** Results of bio reactor metagenome reads classification

| 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 |

327

328   *4.5. Assembly of the bio reactor metagenome before and after reads binning*

329   We assembled reads belonging to different clusters separately, and compared the resulting contigs

330   with the results of the assembly of the entire dataset. The total number of contigs after assembly of

331   the partitioned dataset was comparable to the amount of contigs obtained from the assembly of the

332   entire dataset (Table 5). The same can be said about the total length of contigs and contigs length

333   distributions (see supplementary materials). These results, showing that the database partitioning did

334   not lead to the change of the contigs number or their lengths, can be seen as indirect evidence proving

335   that our $k$-mer based binning of metagenome reads results in species-based clustering.

336   We compared the assembled contigs obtained for the entire and partitioned datasets to the reference

337   genomes of *M. oxyfera*, *K. stuttgartiensis* and *B. sinica*. Even though we could successfully map

338   around 9% of the reads to the reference genomes of *K. stuttgartiensis* and *B. sinica*, we did not get

339  contigs that could be mapped to these genomes. However, the contigs assembled from the entire and

340  partitioned datasets did map to *M. oxyfera* genome. Only 91 out of 196 contigs obtained from the

341  entire dataset assembly could be mapped back to the *M. oxyfera* genome covering 54% of its length.

342  For the assembly of the partitioned dataset, 85 contigs were mapped to the genome of *M. oxyfera* in

343  total, covering 52.65% of its length. The vast majority of those contigs (79, covering 51% of the

344  *M. oxyfera* genome length) derived from the assembly of reads belonging to one cluster. Thus, our

345  dataset partitioning binned the majority of contigs according to their origin.

346

347  **Table 5.** Results of entire and partitioned bio reactor sequencing data assembly and comparison of obtained contigs to

348  the *M.oxyfera* genome.

| Dataset assembled | Entire dataset | Cluster 1 | Cluster 2 | Cluster 3 | Cluster 4 | Cluster 5 | Cluster 6 | Cluster 7 | Cluster 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 | 79 | 4 |
| Contigs mapped on *M.oxyfera* genome | 91 | 0 | 0 | 9 | 1 | 2 | 0 | 71 | 2 |
| Length of mapped contigs, bp | 1 842 182 | 0 | 0 | 132 863 | 11 945 | 21 105 | 0 | 1 497 132 | 17 013 |
| % of *M.oxyfera* genome covered | 54 | 0 | 0 | 1.2 | 0.1 | 0.15 | 0 | 51 | 0.2 |

349

19

350    DISCUSSIONS

351    We described a new approach for efficient, alignment-free binning of metagenomic sequencing reads

352    based on $k$-mer frequencies. Our method successfully classifies reads per organism of origin, for both

353    simulated and real metagenomics data.

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

355    sequencing of a real bio reactor metagenome. The absolute majority of the reads with known origin

356    (*M. oxyfera* or *K. stuttgartiensis/B. sinica*) were clustered together per origin after pairwise

357    comparison of their $k$-mer profiles and subsequent density-based cluster detection. This result was

358    robust, as we observed during the analysis of five subsets of the original PacBio sequencing data with

359    overlapping content. The same experiment demonstrated that each subset provides a similar number

360    of clusters. Reads with unknown origin tended to cluster similarly among different subsets, again

361    confirming the clustering consistency. Although the majority of reads in the analyzed metagenome

362    was of unknown origin, the results can be used to estimate the microbial community complexity for

363    its most abundant inhabitants.

364    The binning of the bio-reactor metagenomics dataset had almost no influence on the results of the

365    metagenome assembly. The number of contigs and their lengths obtained for the entire and partitioned

366    datasets were comparable. This indicates that the $k$-mer based reads binning leads to the organism-

367    based partitioning of metagenomic data. Furthermore, contigs, belonging to the same organism, were

368    automatically grouped together when assembling the dataset subjected to the classification procedure.

369    Thus, our $k$-mer based binning technique can be used to interpret metagenomic assembly results.

370    Performing the binning procedure on an artificially generated PacBio datasets lead to a reads

371    classification per organism, even after adding reads with unknown origin (noise reads). Moreover,

372    increasing the proportion of noise reads leads to a better separation between them and the reads with

373    known origin. This observation supports the central hypothesis of this research, namely that $k$-mer

20

374    distances can be used to cluster reads of the same origin together once those reads provide sufficient

375    coverage of the organisms' genome.

376    The main disadvantages of the current implementation of our method is the limited number of reads

377    (10 000) that can be analyzed. As mentioned before, reads, derived from the same organism, will

378    cluster together, but this is possible only under the condition that the organisms' genome is

379    sufficiently covered. Thus, the described technique is unsuitable for the analysis of metagenomes with

380    a large number of inhabitants or when the inhabitants have large genomes, as 10 000 reads will not

381    be enough to provide sufficient coverage. The depth of the classification that can be performed by the

382    suggested method is still to be discovered.

383    We believe that adapting our metagenomics reads binning technique for larger sets of data and further

384    investigation of its metagenome resolving capacity would allow to expand the current limits of

385    microbiology in the future.

386    CONCLUSIONS

387    In this study we demonstrated the possibility to detect substructures within a single metagenome

388    operating only with the information derived from the sequencing reads. Results obtained for both

389    artificial and real metagenomic data indicated the reads clustering per their known origin. We have

390    shown the robustness of the obtained results by adding different proportions of "noise" reads to the

391    artificially generated metagenomic data and by comparing the results of binning procedure performed

392    on the different subsets of the same real metagenomic dataset. The obtained results are highly

393    important as they establish a principle that might potentially greatly expand the toolkit for the

394    detection and investigation of previously unknow microorganisms.

395    LIST OF ABBREVIATIONS

396    PacBio - Pacific Biosciences

397    NGS - next-generation sequencing;

398    N-DAMO - Nitrite/Nitrate dependent Anaerobic Methane Oxidation

399    Annamox - anaerobic ammonium oxidation

400    WGS - whole-genome shotgun sequencing.

401

402    DECLARATIONS

403    *Ethics approval and consent to participate*

404    Since in this research no human material or clinical records of patients or volunteers were used, this

405    research is out of scope for a medical ethical committee. This information was verified by the Leiden

406    University Medical Center Medical Ethical Committee.

407    *Consent for publication*

408    Not applicable

409    *Availability of data and material*

410    Sequencing reads of bioreactor metagenome were submitted to NCBI under the BioProject number

411    PRJNA487927. Artificial PacBio metagenomic reads with the addition of $0\%, 5\%, 10\%,$ and $15\%$ of

412    real "noise" reads were submitted to NCBI under the BioProject number PRJNA533970.

413    Supplementary materials (Additional file 1) were deposited on Figshare and available for

414    downloading using the following link: https://doi.org/10.6084/m9.figshare.c.4218857.v1.

415    Example of the classification procedure can be found using the following link:

416    https://git.lumc.nl/l.khachatryan/pacbio-meta/blob/master/analysis/real_data/tsne_subset2/analysis_example.ipynb

417    *Competing interests*

418    The authors declare that they have no competing interests

419    *Funding*

424   *Authors' contributions*

425   LK algorithm developing, data acquisition, analysis and interpretation, manuscript drafting; SYA

426   conception, data acquisition and analysis, manuscript editing; RHAMV data acquisition, manuscript

427   editing; JFJL conception, manuscript editing, general supervision.

428   *Acknowledgements*

432    REFERENCES

433    [1] Bikel S, Valdez-Lara A, Cornejo-Granados F, Rico K, Canizales-Quinteros S, Soberon X, et al.

434    Combining metagenomics, metatranscriptomics and viromics to explore novel microbial interactions:

435    towards a systems-level understanding of human microbiome. Computational and structural

436    biotechnology journal 2015;13:390–401.

437    [2] Gosalbes MJ, Abellan JJ, Durban A, Perez-Cobas AE, Latorre A, and Moya A. Metagenomics of

438    human microbiome: beyond 16S rDNA. Clinical Microbiology and Infection 2012;18:47–49.

439    [3] Maccaferri S, Biagi E, and Brigidi P. Metagenomics: key to human gut microbiota. Digestive

440    diseases 2011;29(6):525–530.

441    [4] Martin R, Miquel S, Langella P, and Bermudez-Humaran LG. The role of metagenomics in

442    understanding the human microbiome in health and disease. Virulence 2014;5(3):413–423.

443    [5] Edmonds-Wilson SL, Nurinova NI, Zapka CA, Fierer N, and Wilson M. Review of human hand

444    microbiome research. Journal of dermatological science 2015;80(1):3–12.

445    [6] Blum HE. The human microbiome. Advances in medical sciences 2017;62(2):414–420.

446    [7] Holmes E, Li JV, Marchesi JR, and Nicholson JK. Gut microbiota composition and activity in

447    relation to host metabolic phenotype and disease risk. Cell metabolism 2012;16(5):559–564.

448    [8] Bhatt AP, Redinbo MR, and Bultman SJ. The role of the microbiome in cancer development

449    and therapy. CA: a cancer journal for clinicians 2017;67(4):326–344.

450    [9] Cho I and Blaser MJ. The human microbiome: at the interface of health and disease.

451    Nature Reviews Genetics 2012;13(4):260.

452    [10] Sonnenburg JL and Backhed F. Diet–microbiota interactions as moderators of human

453    metabolism. Nature 2016;535(7610):56.

454    [11] Mullish BH, Marchesi JR, Thursz MR, and Williams HRT. Microbiome manipulation with faecal

455    microbiome transplantation as a therapeutic strategy in clostridium difficile infection. QJM: An

456    International Journal of Medicine 2014;108(5):355–359.

457 [12] Moloney RD, Desbonnet L, Clarke G, Dinan TG, and Cryan JF. The microbiome: stress, health

458 and disease. Mammalian Genome 2014;25(1-2):49–74.

459 [13] Contreras AV, Cocom-Chan B, Hernandez-Montes G, Portillo-Bobadilla T, and Resendis-

460 Antonio O. Host-microbiome interaction and cancer: Potential application in precision medicine.

461 Frontiers in physiology 2016;7:606.

462 [14] He C, Shan Y, and Song W. Targeting gut microbiota as a possible therapy for diabetes. Nutrition

463 Research 2015;35(5):361–367.

464 [15] Marx CJ. Can you sequence ecology? metagenomics of adaptive diversification.

465 PLoS biology 2013;11(2):e1001487.

466 [16] Hiraoka S, Yang C-C, and Iwasaki W. Metagenomics and bioinformatics in microbial ecology:

467 current status and beyond. Microbes and environments 2016;31(3):204–212.

468 [17] Xing M-N, Zhang X-Z, and Huang H. Application of metagenomic techniques in mining

469 enzymes from microbial communities for biofuel synthesis. Biotechnology advances

470 2012;30(4):920–929.

471 [18] Tiwari R, Nain L, Labrou NE, and Shukla P. Bioprospecting of functional cellulases from

472 metagenome for second generation biofuel production: a review. Critical reviews in microbiology

473 2018;44(2):244–257.

474 [19] Sommer MOA, Church GM, and Dantas G. A functional metagenomic approach for expanding

475 the synthetic biology toolbox for biomass conversion. Molecular systems biology 2010;6(1):360.

476 [20] Bokulich NA, Lewis ZT, Boundy-Mills K, Mills DA. A new perspective on microbial landscapes

477 within food production. Current opinion in biotechnology 2016;37:182–189.

478 [21] Trindade M, van Zyl LJ, Navarro-Fernandez J, and Abd Elrazak A. Targeted metagenomics as a

479 tool to tap into marine natural product diversity for the discovery and production of drug candidates.

480 Frontiers in microbiology 2015;6:890.

26

481  [22] Zhang MM, Qiao Y, Ang EL, and Zhao H. Using natural products for drug discovery: the impact

482  of the genomics era. Expert opinion on drug discovery 2017;12(5):475–487.

483  [23] Techtmann SM and Hazen TC. Metagenomic applications in environmental monitoring and

484  bioremediation. Journal of industrial microbiology & biotechnology 2016;43(10):1345–1354.

485  [24] Kunin V, Copeland A, Lapidus A, Mavromatis K, and Hugenholtz P. A bioinformatician's guide

486  to metagenomics. Microbiology and molecular biology reviews 2008;72(4):557–578.

487  [25] Mande SS, Mohammed MH, and Ghosh TS. Classification of metagenomic sequences: methods

488  and challenges. Briefings in bioinformatics 2012;13(6):669–681.

489  [26] Ding X, Cheng F, Cao C, and Sun X. Dectico: an alignment-free supervised metagenomic

490  classification method based on feature extraction and dynamic selection. BMC Bioinformatics 2015

491  16(1):323.

492  [27] Cui H and Zhang X. Alignment-free supervised classification of metagenomes by recursive svm.

493  BMC Genomics 2013;14(1):641.

494  [28] Liao W, Ren J, Wang K, Wang S, Zeng F, Wang Y, and Sun F. Alignment-free transcriptomic

495  and metatranscriptomic comparison using sequencing signatures with variable length Markov chains.

496  Scientific reports 2016;6:37243.

497  [29] Laczny CC, Kiefer C, Galata V, Fehlmann T, Backes C, and Keller A. Busybee web:

498  metagenomic data analysis by bootstrapped supervised binning and annotation. Nucleic acids research

499  2017;45(W1):W171–W179.

500  [30] Wang Y, Hu H, and Li X. Mbmc: An effective Markov chain approach for binning metagenomic

501  reads from environmental shotgun sequencing projects. Omics: a journal of integrative biology

502  2016;20(8):470–479.

503  [31] Kotamarti RM, Hahsler M, Raiford D, McGee M, and Dunham MaH. Analyzing taxonomic

504  classification using extensible Markov models. Bioinformatics 2010;26(18):2235–2241.

505  [32] Seok H-S, Hong W, and Kim J. Estimating the composition of species in metagenomes by

506  clustering of next generation read sequences. Methods 2014;69(3):213–219.

507  [33] Lemos LN, Fulthorpe RR, Triplett EW, and Roesch LFW. Rethinking microbial diversity

508  analysis in the high throughput sequencing era. Journal of microbiological methods 2011;86(1):42–

509  51.

510  [34] Janssen P, Goldovsky L, Kunin V, Darzentas N, and Ouzounis CA. Genome coverage, literally

511  speaking: The challenge of annotating 200 genomes with 4 million publications. EMBO reports 2005;

512  6(5):397–399.

513  [35] Akondi KB and Lakshmi VV. Emerging trends in genomic approaches for microbial

514  bioprospecting. Omics: a journal of integrative biology 201317(2):61–70.

515  [36] Hunter-Cevera JC. The value of microbial diversity. Current Opinion in Microbiology

516  1998;1(3):278–285.

517  [37] Pace NR. Mapping the tree of life: progress and prospects. Microbiology and molecular biology

518  reviews 2009;73(4):565–576.

519  [38] Grattepanche J-D, Santoferrara LF, McManus GB, and Katz LA. Diversity of diversity:

520  conceptual and methodological differences in biodiversity estimates of eukaryotic microbes as

521  compared to bacteria. Trends in microbiology 2014;22(8):432–437.

522  [39] Zinger L, Gobet A, and Pommier T. Two decades of describing the unseen majority of aquatic

523  microbial diversity. Molecular Ecology 2012;21(8):1878–1896.

524  [40] Szalkai B, Scheer I, Nagy K, Vertessy BG, Grolmusz V. The metagenomic telescope. PloS One

525  2014;9(7):e101605.

526  [41] Rosen GL, Polikar R, Caseiro DA, Essinger SD, and Sokhansanj BA. Discovering the unknown:

527  improving detection of novel species and genera from short reads. BioMed Research International

528  2011;2011:495849. doi: 10.1155/2011/495849.

529  [42] Kislyuk A, Bhatnagar S, Dushoff J, and Weitz JS. Unsupervised statistical clustering of

530  environmental shotgun sequences. BMC Bioinformatics 2009;10(1):316.

531  [43] Roumpeka DD, Wallace RJ, Escalettes F, Fotheringham I, and Watson M. A review of

532  bioinformatics tools for bio-prospecting from metagenomic sequence data. Frontiers in genetics

533  2017;8:23.

534  [44] Wu Y-W and Ye Y. A novel abundance-based algorithm for binning metagenomic sequences

535  using l-tuples. Journal of Computational Biology 2011;18(3):523–534.

536  [45] Wang Y, Leung HCM, Yiu S-M, and Chin FYL. Metacluster 4.0: a novel binning algorithm for

537  NGS reads and huge number of species. Journal of Computational Biology 2012;19(2):241–249.

538  [46] Van Lang T, Van Hoai T, et al. A two-phase binning algorithm using l-mer frequency on groups

539  of non-overlapping reads. Algorithms for Molecular Biology 2015;10(1):2.

540  [47] Song K, Ren J, Reinert G, Deng M, Waterman MS, and Sun F. New developments of alignment-

541  free sequence comparison: measures, statistics and next-generation sequencing. Briefings in

542  bioinformatics 2013;15(3):343–353.

543  [48] Girotto S, Pizzi C, and Comin M. Metaprob: accurate metagenomic reads binning based on

544  probabilistic sequence signatures. Bioinformatics 2016;32(17):i567–i575.

545  [49] Anvar SY, Khachatryan L, Vermaat M, van Galen M, Pulyakhina I, Ariyurek Y, Kraaijeveld K,

546  den Dunnen JT, de Knijff P, Ac't Hoen P, et al. Determining the quality and complexity of next-

547  generation sequencing data without a reference genome. Genome Biology 2014;15(12):555.

548  [50] Ono Y, Asai K, and Hamada M. Pbsim: PacBio reads simulator—toward accurate

549  genome assembly. Bioinformatics 2012;29(1):119–121.

550  [51] van Eijk E, Anvar SY, Browne HP, Leung WY, Frank J, Schmitz AM, Roberts AP, and Smits

551  WK. Complete genome sequence of the Clostridium difficile laboratory strain 630d erm reveals

552  differences from strain 630, including translocation of the mobile element ctn 5. BMC Genomics

553  2015;16(1):31.

554  [52] Haroon MF, Hu S, Shi Y, Imelfort M, Keller J, Hugenholtz P, Yuan Z, and Tyson GW. Anaerobic

555  oxidation of methane coupled to nitrate reduction in a novel archaeal lineage. Nature

556  2013;500(7464):567.

557  [53] Chaisson MJ and Tesler G. Mapping single molecule sequencing reads using basic local

558  alignment with successive refinement (blasr): application and theory. BMC Bioinformatics

559  2012;13(1):238.

560  [54] Altschul SF, Gish W, Miller W, Myers EW, and Lipman DJ. Basic local alignment search tool.

561  Journal of molecular biology 1990;215(3):403–410.

562  [55] Chin C-S, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, Dunn C, O'Malley,

563  Figueroa-Balderas RR, Morales-Cruz A, et al. Phased diploid genome assembly with single-molecule

564  real-time sequencing. Nature methods 2016;13(12):1050.

565   [56] Hamada M, Ono Y, Asai K, and Frith MC. Training alignment parameters for arbitrary

566  sequencers with last-train. Bioinformatics 2016;33(6):926–928.

567  [57] van der Maaten L and Hinton G. Visualizing data using t-sne. Journal of machine learning

568  research 2008;9(Nov):2579–2605.

569  [58] Ester M, Kriegel H-P, Sander J, Xu X, et al. A density-based algorithm for discovering clusters

570  in large spatial databases with noise. In Kdd1996 ;96 :226–231.

571  [59] Frank J, Lucker S, Vossen RHAM, Jetten MSM, Hall RJ, Op den Camp HJM, and Anvar SY.

572  Resolving the complete genome of Kuenenia Stuttgartiensis from a membrane bioreactor enrichment

573  using single-molecule real-time sequencing. Scientific reports 2018;8(1):4580.

574

575  ADDITIONAL FILES

576  Additional file 1: Article supplement (PDF 143 kb).

577  Section S1: Threshold for the contig-genome similarity using LAST; Section S2: Detailed results of

578  artificial metagenomes binning.