# GigaScience

## Libra: robust biological inferences of global datasets using scalable k-mer based all-vs-all metagenome comparisons
### --Manuscript Draft--

| | |
|---|---|
| Manuscript Number: | GIGA-D-18-00324 |
| Full Title: | Libra: robust biological inferences of global datasets using scalable k-mer based all-vs-all metagenome comparisons |
| Article Type: | Research |
| Funding Information: | Directorate for Computer and Information Science and Engineering (1640775) — Prof. Bonnie L Hurwitz |

**Abstract:**

Background

Shotgun metagenomics provides powerful insights into microbial community biodiversity and function. Unfortunately, inferences from metagenomic studies are often limited by dataset size and complexity, and are restricted by the availability and completeness of existing databases. De novo comparative metagenomics enables the comparison of metagenomes based on their total genetic content.

Results

We developed a novel tool called Libra that performs all-vs-all comparison of metagenomes based on their k-mer-composition. This tool presents three main innovations: the use of a scalable Apache Hadoop framework enabling massive dataset comparison, the use of complex distance metrics allowing precise clustering of metagenomes based on their k-mer content, and a web-based tool imbedded in iMicrobe (http://imicrobe.us) that uses the CyVerse advanced cyberinfrastructure to promote broad use of the tool by the scientific community.

Conclusions

A comparison of Libra to equivalent tools using both simulated and real metagenomic datasets, ranging from 80 million to 4.2 billion reads, reveals that numerous methods commonly implemented to reduce compute time for large datasets—such as data reduction, read count normalization, and presence/absence distance metrics—greatly diminish the degree of resolution and robustness of large-scale comparative analyses. In contrast, Libra provides scalable high-resolution comparisons using all reads without biases due to differences in abundance and read depth, enabling global-scale analyses to identify microbial signatures linked to biological processes.

| | |
|---|---|
| Corresponding Author: | Bonnie Hurwitz<br><br>UNITED STATES |
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | |
| Corresponding Author's Secondary Institution: | |
| First Author: | Illyoung Choi, MS |
| First Author Secondary Information: | |
| Order of Authors: | Illyoung Choi, MS |
| | Alise J. Ponsero, PhD |
| | Matthew Bomhoff, MS |
| | Ken Youens-Clark, BA |

| | |
|---|---|
| | John H. Hartman, PhD |
| | Bonnie L Hurwitz, PhD |
| **Order of Authors Secondary Information:** | |
| **Additional Information:** | |
| **Question** | **Response** |
| Are you submitting this manuscript to a special series or article collection? | No |
| **Experimental design and statistics**<br><br>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.<br><br>Have you included all the information requested in your manuscript? | Yes |
| **Resources**<br><br>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.<br><br>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist? | Yes |
| **Availability of data and materials**<br><br>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using | Yes |

| a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript.<br><br><br>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist? | |
|---|---|

1 **Title:** Libra: robust biological inferences of global datasets using scalable k-mer based all-vs-all

2 metagenome comparisons

3

4 **Authors:** Illyoung Choi[1], Alise J. Ponsero[2], Matthew Bomhoff[2], Ken Youens-Clark[2], John H.

5 Hartman[1*], and Bonnie L. Hurwitz[2,3*]

6

7 **Affiliations:**

8 [1]Department of Computer Science, University of Arizona, Tucson, Arizona

9 [2]Department of Biosystems Engineering, University of Arizona, Tucson, Arizona

10 [3]BIO5 Institute, University of Arizona, Tucson, Arizona

11 **Corresponding Author:**

12 Bonnie L. Hurwitz bhurwitz@email.arizona.edu

13

14

15

16

17

18

19

**ABSTRACT**

**Background**: Shotgun metagenomics provides powerful insights into microbial community biodiversity and function. Unfortunately, inferences from metagenomic studies are often limited by dataset size and complexity, and are restricted by the availability and completeness of existing databases. *De novo* comparative metagenomics enables the comparison of metagenomes based on their total genetic content.

**Results**: We developed a novel tool called Libra that performs all-vs-all comparison of metagenomes based on their k-mer-composition. This tool presents three main innovations: the use of a scalable Apache Hadoop framework enabling massive dataset comparison, the use of complex distance metrics allowing precise clustering of metagenomes based on their k-mer content, and a web-based tool imbedded in iMicrobe (http://imicrobe.us) that uses the CyVerse advanced cyberinfrastructure to promote broad use of the tool by the scientific community.

**Conclusions**: A comparison of Libra to equivalent tools using both simulated and real metagenomic datasets, ranging from 80 million to 4.2 billion reads, reveals that numerous methods commonly implemented to reduce compute time for large datasets—such as data reduction, read count normalization, and presence/absence distance metrics—greatly diminish the degree of resolution and robustness of large-scale comparative analyses. In contrast, Libra provides scalable high-resolution comparisons using all reads without biases due to differences in abundance and read depth, enabling global-scale analyses to identify microbial signatures linked to biological processes.

**Keywords**: metagenomics, Hadoop, k-mer, distance metrics, clustering

2

43 **INTRODUCTION**

44 Over the last decade, scientists have generated petabytes of genomic data to uncover the role

45 of microbes in dynamic living systems. Yet to understand the underlying biological principles

46 that guide the distribution of microbial communities, massive 'omics datasets need to be

47 compared with environmental factors to find linkages across space and time. One of the

48 greatest challenges in these endeavors has been in documenting and analyzing unexplored

49 genetic diversity in wild microbial communities. For example, fewer than 60% of 40 million non-

50 redundant genes from the Global Ocean Survey (GOS) and the Tara Oceans Expeditions match

51 known proteins in bacteria [1,2]. Other microorganisms such as viruses or pico- eukaryotes that

52 are important to ocean ecosystems are even less well defined (e.g. < 7% of reads from viromes

53 match known proteins [3]). This is largely due to the fact that reference genomes for these

54 organisms do not exist in public data repositories and genome-sequences from metagenomic

55 data await better taxonomic and functional definition. As a result, even advanced tools such as

56 k-mer based classifiers that rapidly assign metagenomic reads to known microbes (Table 1)

57 miss "microbial dark matter" that comprises a significant proportion of metagenomes.

**Table 1.**

| Tool | Area* | Method | Platform | Command line | Parallelized | Scalable** | Web-enabled | Cyber-infrastructure | Cited by | Year |
|---|---|---|---|---|---|---|---|---|---|---|
| **Libra** | **MG** | Pairwise distance calculation | **Hadoop** | **X** | **X** | **X** | **X** | **X** | **current study** | |
| Compareads | MG | Pairwise distance calculation | single server | X | | | | | 35 | 2012 |
| Commet | MG | Pairwise distance calculation | single server | X | | | | | 30 | 2014 |
| Mash | G/MG | Pairwise distance calculation | single server | X | | | | | 157 | 2016 |
| Simka | MG | Pairwise distance calculation | HPC*** | X | X | | | | 18 | 2016 |
| NBC | MG | Taxonomic profiling | singer server | X | | | | | 168 | 2010 |
| Kraken | MG | Taxonomic profiling | singer server | X | | | | | 785 | 2014 |
| FOCUS | MG | Taxonomic profiling | singer server | X | | | X | | 49 | 2014 |
| Clark | MG | Taxonomic profiling | singer server | X | | | | | 176 | 2015 |

| Metaphlan2 | MG | Taxonomic profiling | singer server | X | | | 227 | 2015 |
|---|---|---|---|---|---|---|---|---|
| Metafast | MG | Taxonomic profiling | single server | X | | | 19 | 2016 |
| Centrifuge | MG | Taxonomic profiling | single server | X | | | 78 | 2016 |
| Jellyfish | G/MG | K-mer counting | single server | X | | | 746 | 2011 |
| BioPig | G/MG | K-mer counting | Hadoop | X | X | X | 97 | 2013 |
| Bloomfish | G/MG | K-mer counting | Hadoop | X | X | X | 2 | 2017 |
| Myrna | G | Differential gene expression | Hadoop | X | X | X | 331 | 2010 |
| Eoulsan | G | Differential gene expression | Hadoop | X | X | X | 90 | 2012 |
| Cloud RSD | G | Ortholog detection | Hadoop | X | X | X | 120 | 2010 |
| CloudBLAST | G | Read mapping (ref db) | Hadoop | X | X | X | 362 | 2008 |
| Cloudburst | G | Read mapping (ref genome) | Hadoop | X | X | X | 711 | 2009 |
| Crossbow | G | Variant detection | Hadoop | X | X | X | 501 | 2009 |

\* MG = metagenomics; G = genomics

\** Scalability is defined as reliable distributed high-performance computing framework

\*** High-performance computer

58

59 ***De novo* comparative metagenomics offers a path forward.** In order to examine the

60 complete genomic content, metagenomic samples can be compared using their sequence

61 signature (or frequency of k-mers; Table 1). This approach relies on three core tenets of k-mer-

62 based analytics: (i) closely related organisms share k-mer profiles and cluster together, making

63 taxonomic assignment unnecessary [4,5], (ii) k-mer frequency is correlated with the abundance

64 of an organism [6], and (iii) k-mers of sufficient length can be used to distinguish specific

65 organisms [7]. In 2012, the Compareads [8] method was proposed, followed by Commet [9].

66 Both of these tools compute the number of shared reads between metagenomes using a k-mer-

67 based read similarity measure. The number of shared reads between datasets is then used to

68 compute a Jaccard distance between samples. Given the computational intensity of all-vs-all

69 sequence analysis, several other methods have been employed to reduce the dimensionality of

70 metagenomes and speed up analyses by creating unique k-mer sets and computing the genetic

71 distance between pairs of metagenomes, such as MetaFast [10] and Mash [11]. The fastest of

72 these methods, Mash, indexes samples by unique k-mers to create size-reduced sketches, and

73 compares these sketches using the min-Hash algorithm [12] for computing a genetic distance

4

74 using Jaccard similarity. Yet, the tradeoff for speed is that samples are reduced to a subset of

75 unique k-mers (1k by default) that lack information on k-mer abundance in the samples. Further,

76 given that Mash uses Jaccard similarity only the genetic distance between samples is

77 accounted for (or genetic content in microbial communities) without considering abundance

78 (dominant vs rare organisms in the sample) which is central to microbial ecology and ecosystem

79 processes.

80 Recently, SIMKA [13] was developed to compute a distance matrix between metagenomes by

81 dividing the input datasets into abundance vectors from subsets of k-mers, then rejoining the

82 resulting abundances in a cumulative distance matrix. The methodology can be parallelized to

83 execute the analyses on a high-performance compute cluster (HPC). SIMKA also provides

84 various ecological distance metrics to let the user choose the metric most relevant to their

85 analysis. However, the computational time varies based on the distance metric, where simple

86 distances scale linearly and complex distances metrics scale quadratically as additional

87 samples are added [13]. Moreover, SIMKA normalizes datasets in an all-vs-all comparison by

88 reducing the depth of sequencing for all samples to the least common denominator, therefore

89 decreasing the resolution of the datasets. Lastly, computing k-mer analytics using HPC is

90 subject to reduced fault tolerance for massive datasets.

91 **Scaling sequence analysis using big data analytics via Hadoop.** Hadoop is an attractive

92 platform for performing large-scale sequence analysis because it provides a distributed file

93 system and distributed computation for analyzing massive amounts of data. Hadoop clusters are

94 comprised of commodity servers so that the processing power increases as more computing

95 resources are added. Hadoop also offers a high-level programming abstraction based on

96 MapReduce that greatly simplifies the implementation of new analytical tools. Programmers do

97 not need specialized training in distributed systems and networking to implement distributed

98 programs using Hadoop. Hadoop also provides fault-tolerance by default. When a Hadoop node

99 fails, Hadoop reassigns the failed node's tasks to another node containing a redundant copy of

100 the data those jobs were processing. This differs from HPC where schedulers track failed nodes

101 and either restart the failed computation from the most recent checkpoint, or from the beginning

102 if checkpointing wasn't used. Thus, using a Hadoop infrastructure ensures that computations

103 and data are protected even in the event of hardware failures. These benefits have led to new

104 analytic tools based on Hadoop, making Hadoop a de facto standard in large-scale data

105 analysis. In metagenomics, the development of efficient and inexpensive high-throughput

106 sequencing technologies has led to a rapid increase of the amount of sequence data for

107 studying microbes in diverse environments. However, no Hadoop-enabled comparative

108 metagenomics tools currently exist.

109 Spark [14] is increasingly popular for scientific data analysis [15] because of its outstanding

110 performance provided by fast in-memory processing. Although Libra is currently implemented

111 on Hadoop, Libra can be easily ported to Spark because both Hadoop and Spark have similar

112 interfaces for data processing and partitioning. For example, Resilient Distributed Datasets

113 (RDD) can be partitioned and distributed over a Spark cluster using Libra's k-mer range

114 partitioning. RDDs are memory-resident, allowing Spark to significantly improve the

115 performance of Libra's k-mer counting and distance matrix computation by avoiding slow disk

116 I/O for intermediate data. Nevertheless, we implemented Libra using Hadoop because Spark

117 requires much more RAM than Hadoop, significantly increasing the cost of the cluster.

118 **Existing big data algorithms compare reads to limited genomic reference data**. Recent

119 progress has been made in translating bioinformatics algorithms to big data architectures to

120 overcome scalability issues for genomic but not metagenomic applications (Table 1). Thus far,

121 these algorithms compare large-scale NGS datasets to reference genomic datasets and replace

122 computationally intensive algorithms such as sequence alignment [16], genetic variant detection

123 [17,18], or short read mapping [19–22]. For example, BlastReduce and CloudBurst are parallel

124 sequence mapping tools based on Apache MapReduce [20,21]. These tools, however,

125 implement a query-to-a-reference approach that is inefficient for all-vs-all analyses of reads from

126 metagenomes. Other algorithms such as BioPig [23] and Bloomfish [24] generate an index of

127 sequence data for later partial sequence search and k-mer counting using MapReduce [25].

128 These tools, however, adopt a suffix array approach similar to traditional bioinformatics tools

129 that is inefficient in reading and indexing data on a distributed file system such as Hadoop, thus

130 reducing performance. Moreover, neither tool offers an end-to-end solution for comparing

131 metagenomes consisting of: data distribution on a Hadoop cluster, k-mer indexing and counting,

132 distance computation, and visualization. Finally, none of these tools are enabled in an advanced

133 cyberinfrastructure where users can compute analyses in a simple web-based platform that

134 offers compute, data storage, and analysis tools.

135 **Libra: a tool for scalable all-vs-all sequence analysis in an advanced cyberinfrastructure**

136 Here, we describe a scalable algorithm called Libra that is capable of performing all-vs-all

137 sequence analysis using MapReduce on the Apache Hadoop platform. We demonstrate for the

138 first time that Hadoop can be applied to all-vs-all sequence comparisons of large-scale

139 metagenomic datasets comprised of mixed microbial communities. We present a new distance

140 metric for comparing datasets using Cosine Similarity [34] to consider genetic distance and

141 microbial abundance simultaneously, along with widely accepted distance metrics in biology

142 such as Bray-Curtis [35] and Jensen-Shannon [36]. We validate this new distance metric using

143 simulated metagenomes to show that Libra has exceptional sensitivity in distinguishing complex

144 mixed microbiomes. Next, we show Libra's ability to distinguish metagenomes by both

145 community composition and abundance using 48 samples (16S rRNA and WGS) from the

146 human microbiome project (HMP) across diverse body sites, and compare the results to Mash

147 and SIMKA. Finally, we show that Libra can scale to massive global-scale datasets by

148 examining viral diversity in 43 Tara Ocean Viromes (TOV) from the 2009-2011 Expedition [27]

149 that represent 26 sites containing about 4.2 billion reads. The resulting data demonstrate that

150 Libra provides accurate, efficient, and scalable compute for comparative metagenomics that can

151 be used to discern global patterns in microbial ecology.

7

152 To promote the broad use of the Libra algorithm we developed a web-based tool in iMicrobe

153 (http://imicrobe.us), where users can run Libra using data in their free CyVerse [28,29] account

154 or use datasets that are integrated into the iMicrobe Data Commons. These analyses are

155 fundamental for determining relationships among diverse metagenomes to inform follow-up

156 analyses on microbial-driven biological processes.

157 **DATA DESCRIPTION**

158 **Staggered mock community.** We performed metagenomic shotgun sequencing on a

159 staggered mock community obtained from the Human Microbiome Consortium (HM-277D). The

160 staggered mock community is comprised of genomic DNA from genera commonly found on or

161 within the human body, consisting of 1,000 to 1,000,000,000 16S rRNA gene copies per

162 organism per aliquot. The resulting DNA was subjected to whole genome sequencing as

163 follows. Mixtures were diluted to a final concentration of 1 nanogram/microliter and used to

164 generate whole genome sequencing libraries with the Ion Xpress Plug Fragment Library Kit and

165 manual #MAN0009847, revC (Thermo Fisher Scientific, Waltham, MA, USA). Briefly, 10

166 nanograms of bacterial DNA was sheared using the Ion Shear enzymatic reaction for 12 min

167 and Ion Xpress barcode adapters ligated following end repair. Following barcode ligation,

168 libraries were amplified using the manufacturer's supplied Library Amplification primers and

169 recommended conditions. Amplified libraries were size selected to ~ 200 base pairs using the

170 Invitrogen E-gel Size Select Agarose cassettes as outlined in the Ion Xpress manual and

171 quantitated with the Ion Universal Library quantitation kit. Equimolar amounts of the library were

172 added to an Ion PI Template OT2 200 kit V3. The resulting templated beads were enriched with

173 the Ion OneTouch ES system and quantitated with the Qubit Ion Sphere Quality Control kit (Life

174 Technologies) on a Qubit 3.0 fluorometer (Qubit, NY, NY, USA). Enriched templated beads

175 were loaded onto an Ion PI V2 chip and sequenced according to the manufacturer's protocol

176 using the Ion PI Sequencing 200 kit V3 on a Ion Torrent Proton sequencer. The sequence data

8

177 comprised of ~80 million reads have been deposited to the NCBI Sequence Read Archive under

178 accession SRP115095 under project accession PRJNA397434.

179 **Simulated data derived from the staggered mock community**. The resulting sequence data

180 from the staggered mock community (~80 million reads) were used to develop simulated

181 metagenomes to test the effects of varying read depth, and composition and abundance of

182 organisms in mixed metagenomes. To examine read depth (in terms of raw read counts and file

183 size), we used the known staggered mock community abundance profile to generate an artificial

184 metagenome using GemSim [30] of 2 million reads (454 sequencing) and duplicated the dataset

185 2x, 5x and 10x. We also simulated the effects of sequencing a metagenome more deeply using

186 GemSim [30] to generate simulated metagenomes with 0.5, 1, 5, and 10 million reads based on

187 the relative abundance of organisms in the staggered mock community. Next, we developed

188 four simulated metagenomes to test the effect of changing the dominant organism abundance

189 and genetic composition including: 10 million reads from the staggered mock community (mock

190 1), the mock community with alterations in a few abundant species (mock 2), the mock

191 community with many alterations in abundant species (mock 3), and mock 3 with additional

192 sequences from archaea to further alter the genetic composition (mock 4) as described in

193 Supplemental Table 1. All simulated datasets are available in iMicrobe (http://imicrobe.us).

194 **Human microbiome 16S rRNA gene amplicons and WGS reads**. Human microbiome

195 datasets were downloaded from the NIH Human microbiome project [31] including 48 samples

196 from 5 body sites including: urogenital (posterior fomix), gastrointestinal (stool), oral (buccal

197 mucosa, supragingival plaque, tongue dorsum), airways (anterior nares), and skin

198 (retroauricular crease left and right; Supplemental Table 2). Matched datasets consisting of 16S

199 rRNA reads, WGS reads, and WGS assembled contigs were downloaded from the 16S trimmed

200 dataset and the HMIWGS/HMASM dataset respectively. For the WGS reads dataset, the

201 analysis was run on the paired 1 read file.

202 **Tara ocean viromes**. Tara oceans viromes were downloaded from European Nucleotide

203 Archive (ENA) at EMBL and consisted of 43 viromes from 43 samples at 26 locations across the

204 world's oceans collected during the Tara Oceans (2009-2012) scientific expedition

205 (Supplemental Table 3; [27]). Metadata for the samples was downloaded from PANGAEA [32].

206 These samples were derived from multiple depths including: 16 surface samples (5-6 meters),

207 18 deep chlorophyll maximum samples (DCM; 17-148 meters), and one mesopelagic sample

208 (791 meters). Quality control procedures were applied according to methods described by Brum

209 and colleagues [27].

210 **RESULTS AND DISCUSSION**

211 **Libra computational strategy**. Libra uses Hadoop MapReduce to perform massive all-vs-all

212 sequence comparisons between next-generation sequence (NGS) datasets. Libra is designed

213 to estimate genetic distance accurately without sacrificing performance. Instead, scalable

214 algorithms and efficient resource usage make it feasible to perform all-vs-all comparisons on

215 large datasets.

216 Libra performs all-vs-all distance comparisons using a sweep line algorithm

217 (https://en.wikipedia.org/wiki/Sweep_line_algorithm). Naively, all-vs-all comparisons would

218 require a total of $n \times (n - 1)/2$ comparisons between $n$ samples. Using a sweep line algorithm,

219 Libra can perform these comparisons in a single pass (Supplemental Figure 1). Libra maximizes

220 cluster efficiency using a load balancing algorithm inspired by Terabyte Sort [33] to distribute the

221 workload evenly over the Hadoop cluster. Highly parallelizable inverted index construction and

222 distance matrix computation algorithms enable Libra to scale to any size NGS dataset (often

223 millions of reads), and perform any number of comparisons across datasets, making global

224 ecosystem-level analyses possible.

10

225 **Libra distance calculation.** Libra uses a vector space model to compute the distance between

226 two NGS datasets. In this model each sample is represented by a vector, each dimension of

227 which corresponds to a unique k-mer. Each component of a vector indicates the weight given to

228 the corresponding k-mer in the distance computation. For example, using the frequency (the

229 raw count) of a k-mer as its weight and using 4-mers, the vector <2,4,0,...> indicates that a k-

230 mer 'aaaa' has a weight of two and a k-mer 'aaac' has a weight of four in the sample, etc. The

231 more weight, the more important the k-mer.

232 The distance between two samples can now be measured by comparing their vectors using a

233 distance metric. Libra provides three distance metrics — Cosine Similarity [34], Bray-Curtis [35]

234 and Jensen-Shannon [36]. In this paper, we demonstrate Cosine Similarity as the default

235 distance metric given that it had the shortest runtime for all distances (see Methods).

236 Cosine Similarity determines an estimate of the genetic distance between samples by the angle

237 between the two vectors. The larger the angle, the larger the distance. The cosine is one when

238 the angle is zero (i.e. the vectors are identical except for their magnitude) and less than one

239 otherwise (see Supplemental Methods for a detailed description).

240 The cosine of the angle does not depend on the magnitude (length) of the vectors. This is

241 advantageous in comparing samples with different sizes of samples (or sequencing depth). For

242 example, if there are two samples with the same composition of k-mers but one has k-mers with

243 double the frequency than the other, their vectors will have same angles so that their cosine

244 similarity will one.

245 **Libra implementation.** We implemented Libra on the Hadoop MapReduce platform. This

246 allows Libra to run on any standard Hadoop 2.3 implementation, while taking advantage of the

247 scalability and fault-tolerance features provided by Hadoop. Hadoop allows robust parallel

248 computation over distributed computing resources via its simple programming interface called

11

249 *MapReduce*, while hiding much of the complexity of distributed computing (e.g. node failures).

250 Taking advantage of Hadoop MapReduce, Libra can scale to larger input datasets and more

251 computing resources. Furthermore, many cloud providers such as Amazon and Google offer

252 Hadoop clusters on a pay-as-you-go basis, allowing scientists to scale their Libra computations

253 to match their datasets and budgets.

254 Libra is implemented using three different MapReduce jobs — 1) k-mer histogram construction,

255 2) inverted index construction, and 3) distance matrix computation. Figure 1 shows a workflow

256 of the Libra algorithm.

257 **Figure 1. The Libra Workflow.**

258 Libra consists of three MapReduce jobs (yellow boxes) — 1) k-mer histogram construction, 2)

259 inverted index construction and 3) distance matrix computation. k-mer histograms are first

260 constructed for input samples to balance workloads over the Hadoop cluster during the

261 subsequent jobs. Inverted indices are constructed per a group of samples in parallel by

262 partitioning k-mer ranges. An index chunk is produced from each partition and an inverted index

263 is constructed from multiple index chunks. During the distance matrix computation, partial

264 contributions are computed within a partition and accumulated to produce the final distance

265 matrix.

266 Libra constructs a k-mer histogram of the input samples for load-balancing. A separate Map

267 task is spawned for every data block in the input sample files to calculate the k-mer histogram

268 for each sample. Thus, the k-mer histogram of the input samples is computed in parallel by

269 running multiple Map tasks and a Reduce task that combines their results.

270 Libra performs the inverted index construction in parallel. In the Map phase, a separate Map

271 task is spawned for every data block in the input sample files. Each Map task generates k-mers

272 from the sequences stored in a data block then passes them to the Reduce tasks. In the

273 Reduce phase, the I/O and computation is split by partitioning the k-mer space using the k-mer

274 histograms computed in the first phase (Supplemental Figure 2). A separate Reduce task is

275 spawned for every partition and a custom Partitioner routes the produced k-mers to Reduce

276 tasks by their k-mer ranges. Each Reduce task then counts k-mers it receives and produces an

277 index chunk. As a result, each index chunk is stored as a separate file in the Hadoop MapFile

278 format. The MapFile is well-suited for Libra as it is designed to store key-value pairs in key

279 order, and supports binary search of the keys.

280 In the distance matrix computation, the work is split by partitioning the k-mer space in the

281 beginning of a MapReduce job. The k-mer histogram files for input samples are loaded and

282 merged, and the k-mer space is partitioned according to the k-mer distributions. A separate Map

283 task is spawned for each partition to perform the computation in parallel. As a result, each task

284 produces an output file containing partial contributions to the score matrix. At the end of the job,

285 Libra merges the partial contributions from the files and produces the complete distance matrix.

286 **Advanced cyberinfrastructure for Libra in iMicrobe**. To improve access to Libra we made it

287 available at iMicrobe (https://www.imicrobe.us). A researcher with a CyVerse account can run

288 Libra on iMicrobe by filling-out a simple web form specifying the input files and parameters.

289 Input files are selected from the CyVerse Data Store where they have either been uploaded by

290 the user to their home directory or are part of the iMicrobe Data Commons. When a job is

291 submitted, the user is presented with the status of the job, and on completion the output files

292 and visualization of results. To deploy Libra on iMicrobe, we developed a job dispatch service to

293 automate execution of Libra on a University of Arizona Hadoop cluster.  The service is written in

294 NodeJS and accepts a JSON description of the job inputs and parameters, stages the input files

295 onto the UA Hadoop cluster, executes Libra with the given parameters, and transfers the

296 resulting output files to the user's home directory in the CyVerse Data Store. The service

13

297    provides a RESTful interface that mimics the Agave API Jobs service and is secured using an

298    Agave OAuth2 token.  Source code is located at https://github.com/hurwitzlab/occ-plan-b.

299    **Cosine similarity allows for an accurate and normalized comparison of metagenomes.**

300    Jaccard and Bray-Curtis distance have been extensively used to compare metagenomes based

301    on their sequence signature [10,11,13]. While Mash only computes the Jaccard distance

302    between samples, Simka and Libra implement several classical ecology distances allowing the

303    user to choose the best-suited distance for the considered dataset [13]. Moreover, Libra

304    implements a new distance metric, the cosine similarity. Users can also weight k-mers based on

305    their abundance in Libra (using boolean weighting, natural weighting and logarithmic weighting)

306    to account for differences in microbial community composition and sequencing effort as detailed

307    below.

308    We tested these effects by varying: (1) the size of the datasets, (2) depth of sequencing, (3) the

309    abundance of dominant microbes in the community, and (4) genetic composition of the

310    community by adding in an entirely new organism (in our case we added archaea). We

311    constructed simulated metagenomes and compared Libra's distance based on the cosine

312    similarity against those from Mash and SIMKA. Simulated datasets were derived from genomic

313    DNA from a staggered mock community of bacteria obtained from the human microbiome

314    consortium and sequenced deeply using the Ion Torrent sequencing platform (80 million reads,

315    see Methods).

316    First, we examined the effect of the size of the dataset by using GemSim [30] to obtain  a

317    simulated metagenome composed of 1 million reads from the mock community and duplicating

318    that dataset 2x and 10x. Overall, we found that altering the size of the metagenome (by

319    duplicating the data) had no effect on the distance between metagenomes for Mash, SIMKA, or

320    Libra. In each case the distance of the duplicated datasets to the 1x mock community was less

321    than 0.0001 (data not shown).

14

322 Because metagenomes don't scale exactly with size and instead have an increasing

323 representation of low-abundance organisms, we created a second simulated dataset from the

324 mock community using GemSim [30] 0.5, 1, 5, and 10 million reads (454 sequencing) to mimic

325 the effect of sequencing more deeply. Given the abundance of organisms in the mock

326 community, the 0.5 M read dataset is mainly comprised of dominant species. With increased

327 sequencing depth (1, 5, and 10 M reads) additional species are added relative to their

328 abundance in the mock community. Overall, sequencing depth has little effect on the distance

329 between samples in Mash and Libra (natural weighting), whereas SIMKA shows no changes

330 between samples when using Jaccard and Bray-Curtis distances (Figure 2A). Indeed, SIMKA

331 normalization is implemented as follows: the smallest sample from the dataset is determined

332 and its number of sequences is used to compare the samples (in this experiment, all mock

333 communities were compared based on the first 0.5 million reads). These results suggest that

334 Libra (natural weighting) and Mash are appropriate for comparing datasets at different

335 sequencing depths, whereas using SIMKA could lead to undesired effects.

336 **Figure 2. Analysis of artificial metagenomes using Mash, SIMKA and Libra.**

337 A. Distance to staggered mock community artificial metagenome composed of 10 million

338 reads (mock1 10M), for artificial metagenomes of same community sequenced at

339 various depth. Artificial metagenomes were obtained using GemSim and the known

340 abundance profile of the staggered mock community (see Supplemental Table 1). In

341 order to mimic various sequencing depth, the artificial metagenomes were generated at

342 0.5, 1, 5 or 10 million reads (noted mock1 0.5M; mock1 1M; mock1 5M; mock1V2 10M).

343 The distances between the 4 artificial metagenomes and a 10 million read artificial

344 metagenome (mock1 10M) were computing using Mash, SIMKA (Jaccard and Bray-

345 curtis distance) and Libra (natural weighting).

346 B. Distance to staggered mock community artificial metagenome (mock 1), for artificial

15

347    metagenomes from increasingly distant communities. The mock 1 relies on the known

348    abundance profile from the staggered mock community. The mock 2 community profile

349    was obtained by randomly inverting 3 species abundance from mock 1 profile. The mock

350    3 profile was obtained by randomly inverting 2 species abundances from mock 2 profile.

351    Finally, mock 4 profile was obtained by adding high abundance archeal genomes not

352    present in any the other mock communities. Artificial metagenomes were generated

353    using GemSim at 10 million reads. The distance between the mock 1 community to

354    mock 2, mock 3, mock 4 and a replicate community (mock1 V2) was computed using

355    Mash, SIMKA (Jaccard and Bray-curtis distance) and LIBRA (cosine distance, natural

356    and logarithmic weighting).

357  In addition to natural variation in population-level abundances, artifacts from sequencing can

358  result in high-abundance k-mers. Libra allows users to select the optimal methodology for

359  weighting high abundance k-mers in their datasets including boolean, natural, and logarithmic.

360  These options for weighting k-mers are important for different biological scenarios as described

361  below and shown in simulated datasets. To examine the effect of weighting, we compared and

362  contrasted the natural and logarithmic weight in Libra, with other distances obtained from Mash

363  and SIMKA (Jaccard and Bray-Curtis). We also examined the effect of adding an entirely new

364  species by spiking a simulated dataset with sequences derived from archaea (that were not

365  present in the mock community). The simulated datasets were comprised of the staggered

366  mock community (mock 1), the mock community with alterations in a few abundant species

367  (mock 2), the mock community with many alterations in abundant species (mock 3), and mock 3

368  with additional sequences from archaea to alter the genetic composition of the community

369  (mock 4; see Supplemental Table 1). The resulting data showed that Libra (logarithmic

370  weighting) shows a stepwise increase in distance among the mock communities (Figure 2B).

371  This suggests that logarithmic weighting in Libra allows for a comparison of distantly related

372  microbial communities. Mash also shows a stepwise distance between communities, but is

16

373 compressed relative to Libra, making differences less distinct. SIMKA (Bray-Curtis and Jaccard)

374 and Libra (cosine distance, natural weighting) reach the maximum difference between mock

375 communities 3 and 4 (Figure 2B). This indicates that these distances are more appropriate

376 when comparing metagenomes with small fluctuations in the community (e.g., data from a time-

377 series analysis), whereas Libra (cosine distance, logarithmic weighting) can be used to

378 distinguish metagenomes that vary in both genetic composition and abundance over a wide-

379 range of species diversity by dampening the effect of high-abundance k-mers. Because of this

380 important difference, we used the cosine distance with the logarithmic weighting in all

381 subsequent analyses. Cosine distance also provided the fastest computation for complex

382 distance metrics (see Methods).

383 **Libra accurately profiles differences in bacterial diversity and abundance in amplicon**

384 **and WGS datasets from the human microbiome.**

385 Microbial diversity is traditionally assessed using two methods: the 16S rRNA gene to classify

386 bacterial and archaeal groups at the genus to species level, or whole genome shotgun

387 sequencing (WGS) for finer taxonomic classification at the species or subspecies level. Further,

388 WGS datasets provide additional information on functional differences between metagenomes.

389 Here we compare and contrast the effect of different algorithmic approaches (Mash vs Libra vs

390 SIMKA), distance metric (Libra vs SIMKA), data type (16S rRNA vs WGS), and sequence type

391 (WGS reads vs assembled contigs) in analyzing data from 48 samples across 8 body sites from

392 the Human Microbiome Project. Specifically, we examine matched datasets (16S rRNA reads,

393 WGS reads, and WGS assembled contigs) classified as urogenital (posterior fomix),

394 gastrointestinal (stool), oral (buccal mucosa, supragingival plaque, tongue dorsum), airways

395 (anterior nares), and skin (retroauricular crease left and right; Supplemental Table 2).

396 Because the HMP datasets represent microbial communities, abundant bacteria will have more

397 total read counts than rare bacteria in the samples. Thus, each sample can vary by both taxonomic

17

398 composition (the genetic content of taxa in a sample) and abundance (the relative proportion of

399 those taxa in the samples). Importantly, the 16S rRNA amplicon dataset is useful in showing how

400 well each algorithm performs in detecting and quantifying small-scale variation for single a gene at

401 the genus-level, whereas the WGS dataset demonstrates the effect of including the complete

402 genetic content and abundance of organisms at the species-level in a community [37]. Also, we

403 examine differences in each algorithm when read abundance is excluded using assembled contigs

404 that only represent the genetic composition of the community.

405 Using the 16S rRNA reads, both Mash and Libra clustered samples by broad categories but not

406 individual body-sites (Figure 3A and B). Similar to what is described in previous work [13], samples

407 from the airways and skin co-cluster, whereas other categories including urogenital,

408 gastrointestinal, and oral are distinct [13]. These results indicate that limited variation in the 16S

409 rRNA gene may only allow for clustering for broad categories. Further, the Mash algorithm shows

410 lower overall resolution (Figure 3A) as compared to Libra (Figure 3B). Indeed, amplicon

411 sequencing analysis is not an intended use of Mash, given that it reduces the dimensionality of the

412 data by looking at presence/absence of unique k-mers, whereas Libra examines the complete

413 dataset accounting for both composition in organisms and their abundance. In contrast, SIMKA

414 (Jaccard-ab and Bray-Curtis) failed to cluster samples by broad categories: some skin samples are

415 found associated with stool and formix samples (Figure 3C and D). Moreover, SIMKA Jaccard-ab

416 fails to cluster the mouth samples together (Figure 3C). This result suggests that applying SIMKA

417 and these well-used distance metrics are not appropriate for these datasets.

418 **Figure 3. Clustering of HMP 16S rRNA datasets using Mash, Libra and SIMKA.**

419 48 Human metagenomic samples from the HMP projects clustered by Mash (A), Libra (B) or

420 SIMKA using Jaccard-ab (C) and Bray-Curtis distances (D) from 16s sequencing runs. The

421 samples were clustered using Ward's method on their distance scores. Heat maps illustrate the

18

422     pairwise dissimilarity between samples, scaled between 0 (green) and 1 (red). A key below the

423     heatmap colors the samples by body sites.

424     When using WGS reads, both Mash and Libra show enhanced clustering by body-site (Figure 4A

425     and B), however Mash shows decreased resolution (Figure 4A) as compared to Libra (Figure 4B).

426     Again, these differences reflect the effect of using all of the read data (Libra) rather than a subset

427     (Mash). Importantly, the Libra algorithm also depends on read abundance that provides increased

428     resolution for interpersonal variation as seen in skin samples (Figure 4B). Similar to the 16S rRNA

429     datasets, SIMKA (Jaccard-ab and Bray-Curtis) failed to cluster the samples by body site, where

430     some skin and stool samples cluster with formix samples (Figure 4C and D). Similarly, SIMKA

431     Jaccard-ab also fails to cluster the mouth samples together (Figure 4C). Overall SIMKA shows an

432     enhanced clustering by body-site using WGS data compared to the 16S rRNA data using these

433     distance metrics, however the clustering is still not accurate.

434     **Figure 4. Clustering of WGS samples using Mash, and Libra and SIMKA.**

435     48 Human metagenomic samples from the HMP projects clustered by Mash (A), Libra (B) or

436     Simka using Jaccard-ab (C) and Bray-Curtis distances (D) from whole genome shotgun

437     sequencing runs. The samples were clustered using Ward's method on their distance scores.

438     Heat maps illustrate the pairwise dissimilarity between samples, scaled between 0 (green) and

439     1 (red). A key below the heatmap colors the samples by body sites.

440     When abundance is taken out of the equation by using assembled contigs (Supplemental Figure 3)

441     Mash performs well in clustering distinct body sites whereas Libra shows discrepancies and less

442     overall resolution. Thus, Libra requires reads rather than contigs to perform accurately and obtain

443     high-resolution clustering (Figure 4). SIMKA (Jaccard-ab and Bray-Curtis) was not able to

444     distinguish any assembled datasets and scored all sample-to-sample distances to the maximum,

445     even considering presence-absence distance metric proposed by SIMKA (data not shown). This

19

446 phenomenon may be explained by the normalization method used by SIMKA, which does not

447 provide enough data to compare the samples when normalized by the smallest number of contigs

448 (in our dataset 69).

**449 Libra allows for ecosystem-scale analysis: clustering the Tara ocean viromes to unravel**

**450 global patterns.**

451 To demonstrate the scale and performance of the Libra algorithm, we analyzed 43 Tara Ocean

452 Viromes (TOV) from the 2009-2011 Expedition [27] representing 26 sites, 43 samples, and 4.2

453 billion reads from the global ocean (see methods). Phages (viruses that infect bacteria) are

454 abundant in the ocean [38] and can significantly impact environmental processes through host

455 mortality, horizontal gene transfer, and host-gene expression. Yet, how phages change over

456 space and time in the global ocean and with environmental fluxes is just beginning to be

457 explored. The primary challenge is the majority of reads in viromes (often > 90%) do not match

458 known proteins or viral genomes [3] and no conserved genes like the bacterial 16S rRNA gene

459 exist to differentiate populations. To examine known and unknown viruses simultaneously,

460 viromes are best compared using sequence signatures to identify common viral populations.

461 Two approaches exist to cluster viromes based on sequence composition. The first approach

462 uses protein clustering to examine functional diversity in viromes between sites [3,27,39].

463 Protein clustering, however, depends on accurate assembly and gene finding that can be

464 problematic in fragmented and genetically diverse viromes [40]. Further, assemblies from

465 viromes often only include a fraction of the total reads (e.g., only ⅓ in TOV [27]). To examine

466 global viral diversity in the ocean using all of the reads we examined TOV using Libra. The

467 complete pairwise analysis of ~4.2 billion reads in the TOV dataset [27] finished in 18 hours

468 using a 10-node Hadoop cluster (see Methods and Table 2). Importantly, Libra exhibits

469 remarkable performance in computing similarity scores, wherein k-mer matches for all TOV

470 completed within 1.5 hours (Table 2). This step usually represents the largest computational

20

471 bottleneck for bioinformatics tools that compute pairwise distances between sequence pairs for

472 applications such as hierarchical sequence clustering [41–44].

473 *Table 2. Execution times for the Libra based on the Tara Ocean Virome (TOV) dataset.*

| Stage | Execution Time |
|---|---|
| Preprocessing (k-mer histogram construction + Inverted index construction) | 16:32:55 |
| Distance matrix computation | 1:24:27 |
| Total | 17:57:22 |

474

475 Overall, we found that viral populations in the ocean are largely structured by temperature in

476 four gradients (Figure 5) similar to their bacterial hosts [2]. Interestingly, samples from different

477 Longhurst Provinces but the same temperature gradient cluster together. Also, water samples

478 from the surface (SUR) and deep chlorophyll maximum (DCM) at the same station, cluster more

479 closely together than samples from the same depth at nearby sites (Figure 5). Also noteworthy,

480 samples that were derived from extremely cold environments (noted as C0 in Figure 5) lacked

481 similarity to all other samples (at a 30% similarity score), indicating distinctly different viral

482 populations. These samples include a mesotrophic sample that have previously been shown to

483 have distinctly different viral populations than surface ocean samples [45]. Taken together,

484 these data indicate that viral populations are structured globally by temperature, and at finer

485 resolution by station (for surface and DCM samples) indicating that micronutrients and local

486 conditions play an important role in defining viral populations.

487 **Figure 5. Visualizing the genetic distance among marine viral communities using Libra.**

21

488 Distance computed from 43 TOV from the 2009-2012 Tara Oceans Expedition. Lines (edges)

489 between samples represent the similarity and are colored and thickened accordingly. Lines with

490 insignificant similarity (less than 30%) are removed. Each of the sample names are color coded

491 by Longhurst Province. Inner circles show temperature ranges. Sample names show the

492 temperature range, station, and depth as indicated on the legend.

493

494 **INNOVATIONS**

495 Scientific collaboration is increasingly data driven given large-scale next generation sequencing

496 datasets. It is now possible to generate, aggregate, archive, and share datasets that are

497 terabytes and even petabytes in size. Scalability of a system is becoming a vital feature that

498 decides feasibility of massive 'omic's analyses. In particular, this is important for metagenomics

499 where patterns in global ecology can only be discerned by comparing the sequence signatures

500 of microbial communities from massive 'omics datasets, given that most microbial genomes

501 have not been defined. Current algorithms to perform these tasks run on local workstations or

502 high-performance computing architectures that cannot scale. Libra presents three main

503 innovations: the use of a scalable Apache Hadoop framework enabling massive dataset

504 comparison, the use of sophisticated distance metrics allowing high accuracy and clustering of

505 the metagenomes based on their k-mer content, and a web-based tool imbedded in the

506 CyVerse advanced cyberinfrastructure through iMicrobe (http://imicrobe.us) for broader use of

507 the tool in the scientific community. The work described here is the first step in implementing a

508 cloud-based resource for comparative metagenomics that can be broadly used by scientists to

509 analyze large-scale shared data resources. Moreover, the code can be ported to any

510 MapReduce cluster (e.g., Wrangler at TACC, Amazon EMR or private Hadoop clusters). This

511 computing paradigm is consistent with recent efforts to increase the accessibility of big datasets

512 in the cloud, such as the Pan Cancer Analyses of Whole Genomes Project [46].

**METHODS**

513

**Scalability benchmarking for Libra.** We used synthetic datasets for a scalability benchmark. 514

The synthesized datasets consisted of different number of samples, each of which is 10 billion 515

bytes (approximately 9.3 GB). We took samples that are larger than 10 billion bytes from Tara 516

ocean virome dataset and truncated each of them to approximately 10 billion bytes in size while 517

respecting read boundaries. We varied the number of samples to show the scalability of Libra. 518

We used four datasets consisting of 10, 20, 30 and 40 samples in the benchmark. Total sizes of 519

the datasets are 93GB, 186GB, 279GB and 372GB respectively. Each experiment was run 520

three times, and an average of the three runs reported (Supplemental Table 4). 521

**Figure 6**. Scalability testing for Libra. Four datasets consisting of 10, 20, 30 and 40 samples 522

with total sizes of 93GB, 186GB, 279GB and 372GB, respectively. Runtime of Libra increased 523

linearly with increased input volume and number of input samples. The linear increase of 524

runtime shows that Libra efficiently handles increased volume of input and efficiently computes 525

distances between all sample pairs while the number of sample pairs increases quadratically. 526

**Benchmarking runtimes of different distance metrics in Libra.** We used the same synthetic 527

dataset with 40 samples (372GB in total) in the scalability benchmarking. We varied the 528

distance metrics and measured the runtimes of Libra. Because all distance metrics share the 529

same index, we reused the index constructed during the scalability benchmarking, thus, 530

runtimes of the inverted index construction for the different metrics are the same. Each 531

experiment was run three times, and an average of the three runs reported (Supplemental Table 532

4). 533

**Figure 7**. Runtimes of three different distance metrics (Cosine Similarity, Bray-Curtis and 534

Jensen-Shannon) in Libra with 40 samples of input (372GB in total). Differences in runtimes are 535

mainly due to different computational workload of distance metrics. For example, Jensen- 536

23

537 Shannon requires more multiplications and divisions in nested loops than cosine similarity,

538 incurring more computational workload. Yet, distance matrix computation with Jensen-Shannon

539 took only 12.64% of total runtime.

540 **Experimental Environment Description:**

541 **Mash and SIMKA configurations.** Mash v1.1 was run on the metagenomic datasets with the

542 following parameters: -r –s 10000 –m 2 [19]. The analysis of assemblies was run without the

543 parameter "-r", used for short sequences.

544 SIMKA v1.3.2 was run on the metagenomic datasets with the following parameters: -

545 abundance-min 2 -max-reads [MINCOUNT] -simple-dist -complex-dist, where [MINCOUNT] is

546 the smallest sequence count across the analyzed samples.

547 **Hadoop cluster configuration**. The Libra experiments described in the paper were performed

548 on a Hadoop cluster consisting of 10 physical nodes (9 MapReduce worker nodes). Each node

549 contains 12 CPUs and 128 GB of RAM, and is configured to run a maximum of 7 YARN

550 containers simultaneously with 10 GB of RAM per container. The remaining system resources

551 are reserved for the operating system and other Hadoop services such as Hive or Hbase.

560

561 *Competing interests*: The authors declare no competing interests.

**Availability and Implementation**:

Project name: Libra
Project home page: http://github.com/iychoi/libra
Operating system(s): Hadoop 2.3 or higher
Programming language: Java
Other requirements: Java 1.7 or higher
License: Apache License Version 2.0
Any restrictions to use by non-academics: No restriction
Libra web-based App is in iMicrobe under Apps (http://imicrobe.us); Code to implement the
Libra web-based App is in Github (https://github.com/hurwitzlab/occ-plan-b).

**REFERENCES**

1. Yooseph S, Sutton G, Rusch DB, Halpern AL, Williamson SJ, Remington K, et al. The Sorcerer II Global Ocean Sampling expedition: expanding the universe of protein families. PLoS Biol. 2007;5:e16.

2. Sunagawa S, Coelho LP, Chaffron S, Kultima JR, Labadie K, Salazar G, et al. Structure and function of the global ocean microbiome. Science. 2015;348.

3. Hurwitz BL, Sullivan MB. The Pacific Ocean Virome (POV): a marine viral metagenomic dataset and associated protein clusters for quantitative viral ecology. PLoS One. 2013;8:e57355.

4. Dubinkina VB, Ischenko DS, Ulyantsev VI, Tyakht AV, Alexeev DG. Assessment of k-mer spectrum applicability for metagenomic dissimilarity analysis. BMC Bioinformatics. 2016;17:38.

5. Teeling H, Waldmann J, Lombardot T, Bauer M, Glockner FO. TETRA: a web-service and a stand-alone program for the analysis and comparison of tetranucleotide usage patterns in DNA sequences. BMC Bioinformatics. 2004;5:163.

6. Wu Y-W, Ye Y. A novel abundance-based algorithm for binning metagenomic sequences using l-tuples. J Comput Biol. online.liebertpub.com; 2011;18:523–34.

7. Fofanov Y, Luo Y, Katili C, Wang J, Belosludtsev Y, Powdrill T, et al. How independent are the appearances of n-mers in different genomes? Bioinformatics. 2004;20:2421–8.

8. Maillet N, Lemaitre C, Chikhi R, Lavenier D, Peterlongo P. Compareads: comparing huge metagenomic experiments. BMC Bioinformatics. 2012;13 Suppl 19:S10.

9. Maillet N, Collet G, Vannier T, Lavenier D, Peterlongo P. Commet: Comparing and combining multiple metagenomic datasets. 2014 IEEE International Conference on Bioinformatics and Biomedicine (BIBM). 2014. p. 94–8.

10. Ulyantsev VI, Kazakov SV, Dubinkina VB, Tyakht AV, Alexeev DG. MetaFast: fast reference-free graph-based comparison of shotgun metagenomic data. Bioinformatics. 2016;32:2760–7.

11. Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, et al. Mash: fast genome and metagenome distance estimation using MinHash. Genome Biol. 2016;17:132.

12. Chum O, Philbin J, Zisserman A. Near Duplicate Image Detection: min-Hash and tf-idf Weighting. BMVC. 2008; Available from: http://www.bmva.org/bmvc/2008/papers/119.pdf

13. Benoit G, Peterlongo P, Mariadassou M, Drezen E, Schbath S, Lavenier D, et al. Multiple comparative metagenomics using multiset k-mer counting. PeerJ Comput Sci. 2016;2:e94.

14. Zaharia M, Chowdhury M, Franklin MJ, Shenker S, Stoica I. Spark: Cluster computing with working sets. HotCloud. static.usenix.org; 2010;10:95.

15. Guo R, Zhao Y, Zou Q, Fang X, Peng S. Bioinformatics applications on Apache Spark. Gigascience. 2018.

609 16. Kolker N, Higdon R, Broomall W, Stanberry L, Welch D, Lu W, et al. Classifying proteins into
610 functional groups based on all-versus-all BLAST of 10 million proteins. OMICS.
611 online.liebertpub.com; 2011;15:513–21.
612 17. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The
613 Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA
614 sequencing data. Genome Res. genome.cshlp.org; 2010;20:1297–303.
615 18. Langmead B, Schatz MC, Lin J, Pop M, Salzberg SL. Searching for SNPs with cloud
616 computing. Genome Biol. 2009;10:R134.
617 19. Nguyen T, Shi W, Ruden D. CloudAligner: A fast and full-featured MapReduce based tool
618 for sequence mapping. BMC Res Notes. 2011;4:171.
619 20. Schatz MC. BlastReduce: high performance short read mapping with MapReduce.
620 University of Maryland, Available from:
621 https://www.cs.umd.edu/sites/default/files/scholarly_papers/MichaelSchatz_1.pdf
622 21. Schatz MC. CloudBurst: highly sensitive read mapping with MapReduce. Bioinformatics.
623 2009;25:1363–9.
624 22. Pandey RV, Schlötterer C. DistMap: a toolkit for distributed short read mapping on a
625 Hadoop cluster. PLoS One. 2013;8:e72614.
626 23. Nordberg H, Bhatia K, Wang K, Wang Z. BioPig: a Hadoop-based analytic toolkit for large-
627 scale sequence data. Bioinformatics. 2013;29:3014–9.
628 24. Gao T, Guo Y, Wei Y, Wang B, Lu Y, Cicotti P, et al. Bloomfish: A Highly Scalable
629 Distributed K-mer Counting Framework. 2017 IEEE 23rd International Conference on Parallel
630 and Distributed Systems (ICPADS). 2017. p. 170–9.
631 25. Menon RK, Bhat GP, Schatz MC. Rapid Parallel Genome Indexing with MapReduce.
632 Proceedings of the Second International Workshop on MapReduce and Its Applications. New
633 York, NY, USA: ACM; 2011. p. 51–8.
634 26. Salton G, Wong A, Yang CS. A Vector Space Model for Automatic Indexing. Commun ACM.
635 New York, NY, USA: ACM; 1975;18:613–20.
636 27. Brum JR, Ignacio-Espinoza JC, Roux S, Doulcier G, Acinas SG, Alberti A, et al. Patterns
637 and ecological drivers of ocean viral communities. Science. 2015;348.
638 28. Goff SA, Vaughn M, McKay S, Lyons E, Stapleton AE, Gessler D, et al. The iPlant
639 Collaborative: Cyberinfrastructure for Plant Biology. Front Plant Sci. 2011;2:34.
640 29. Devisetty UK, Kennedy K, Sarando P, Merchant N, Lyons E. Bringing your tools to CyVerse
641 Discovery Environment using Docker. F1000Res. 2016;5:1442.
642 30. McElroy KE, Luciani F, Thomas T. GemSIM: general, error-model based simulator of next-
643 generation sequencing data. BMC Genomics. 2012;13:74.
644 31. Human Microbiome Project Consortium. Structure, function and diversity of the healthy
645 human microbiome. Nature. 2012;486:207–14.
646 32. Diepenbroek M, Grobe H, Reinke M, Schindler U, Schlitzer R, Sieger R, et al. PANGAEA—
647 an information system for environmental sciences. Comput Geosci. 2002;28:1201–10.
648 33. O'Malley O. Terabyte sort on apache hadoop. Yahoo, available online at:
649 http://sortbenchmark org/Yahoo-Hadoop pdf,(May). Citeseer; 2008;1–3.
650 34. Huang A. Similarity measures for text document clustering. Proceedings of the sixth new
651 zealand computer science research student conference (NZCSRSC2008), Christchurch, New
652 Zealand. 2008. p. 49–56.
653 35. Michie MG. Use of the Bray-Curtis similarity measure in cluster analysis of foraminiferal
654 data. Math Geol. 1982;14:661–7.
655 36. Lin J. Divergence measures based on the Shannon entropy. IEEE Trans Inf Theory.
656 1991;37:145–51.
657 37. Watts GS, Youens-Clark K, Slepian MJ, Wolk DM, Oshiro MM, Metzger GS, et al. 16S rRNA
658 gene sequencing on a benchtop sequencer: accuracy for identification of clinically important
659 bacteria. J Appl Microbiol. 2017;123:1584–96.

660 38. Bergh O, Borsheim KY, Bratbak G, Heldal M. High abundance of viruses found in aquatic
661 environments. Nature. 1989;340:467–8.
662 39. Hurwitz BL, Brum JR, Sullivan MB. Depth-stratified functional and taxonomic niche
663 specialization in the "core"and "flexible"Pacific Ocean Virome. ISME J. 2014.
664 40. Minot S, Wu GD, Lewis JD, Bushman FD. Conservation of gene cassettes among diverse
665 viruses of the human gut. PLoS One. 2012;7:e42342.
666 41. Sun Y, Cai Y, Huse SM, Knight R, Farmerie WG, Wang X, et al. A large-scale benchmark
667 study of existing algorithms for taxonomy-independent microbial community analysis. Brief
668 Bioinformatics. 2012;13:107–21.
669 42. Edgar RC. Search and clustering orders of magnitude faster than BLAST. Bioinformatics.
670 2010;26:2460–1.
671 43. Niu BF, Fu LM, Sun SL, Li WZ. Artificial and natural duplicates in pyrosequencing reads of
672 metagenomic data. BMC Bioinformatics. 2010;11:187.
673 44. Cai Y, Sun Y. ESPRIT-Tree: hierarchical clustering analysis of millions of 16S rRNA
674 pyrosequences in quasilinear computational time. Nucleic Acids Res. 2011;39:e95.
675 45. Hurwitz BL, Brum JR, Sullivan MB. Depth-stratified functional and taxonomic niche
676 specialization in the "core" and "flexible" Pacific Ocean Virome. ISME J. 2015;9:472–84.
677 46. Weinstein JN, Collisson EA, Mills GB, Shaw KRM, Ozenberger BA, Ellrott K, et al. The
678 cancer genome atlas pan-cancer analysis project. Nat Genet. 2013;45:1113–20.
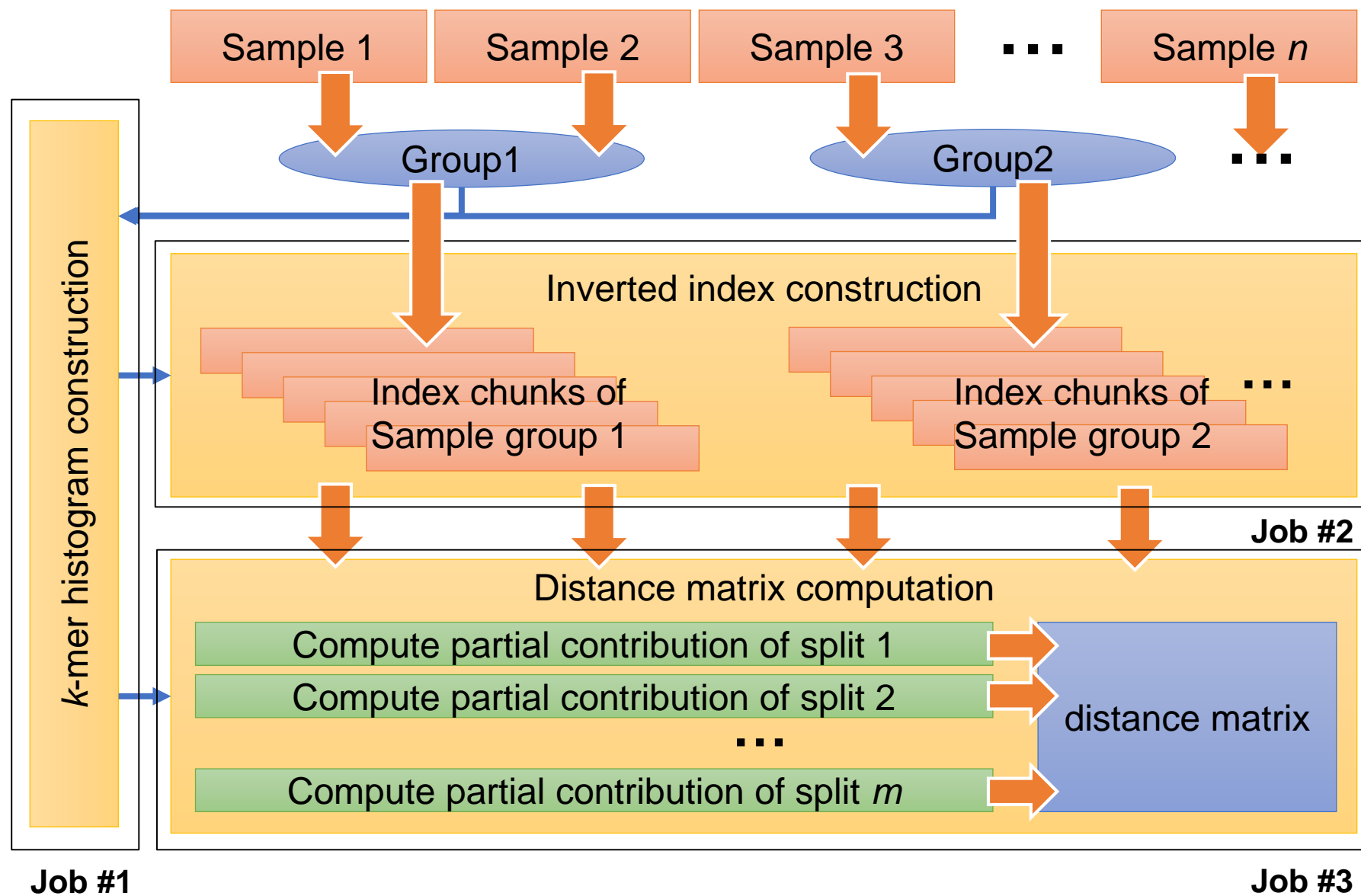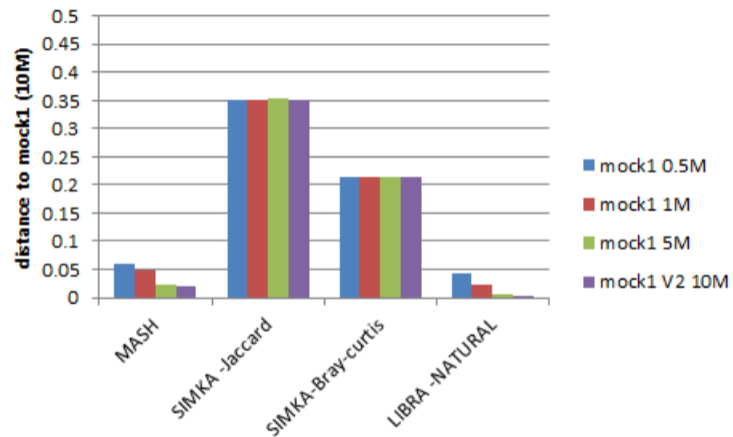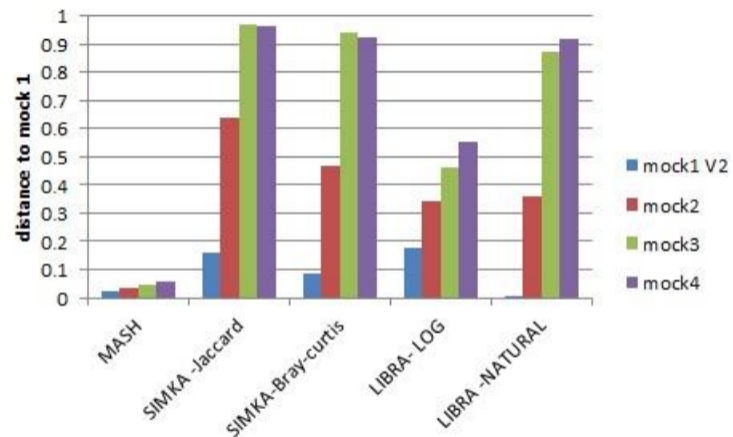679 47. Wrangler - Texas Advanced Computing Center. [cited 2017 Dec 20]. Available from:
680 https://www.tacc.utexas.edu/systems/wrangler

Figure 1

Figure 2

a



b

Figure 3

# a - MASH

# b – LIBRA, log weighting

# c- SIMKA, abundance Jaccard

# d- SIMKA, abundance Bray-Curtis



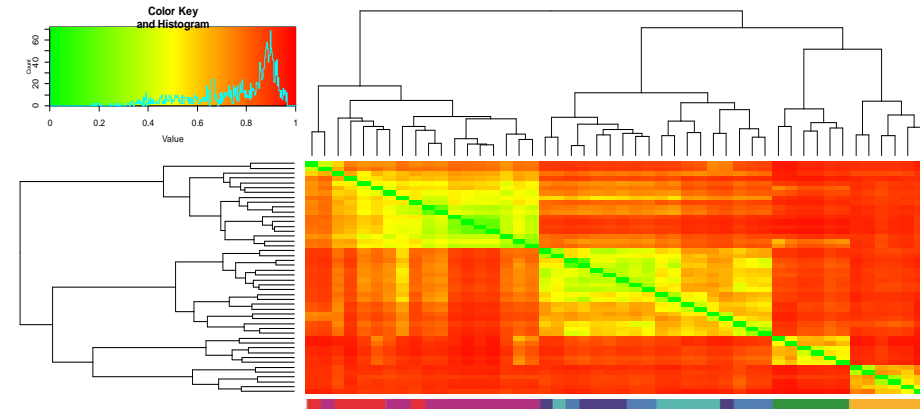Posterior Formix — Stool — Buccal mucosa — Supragingival plaque — Tongue dorsum — Anterior Nares — Retro-auricular crease, left and right
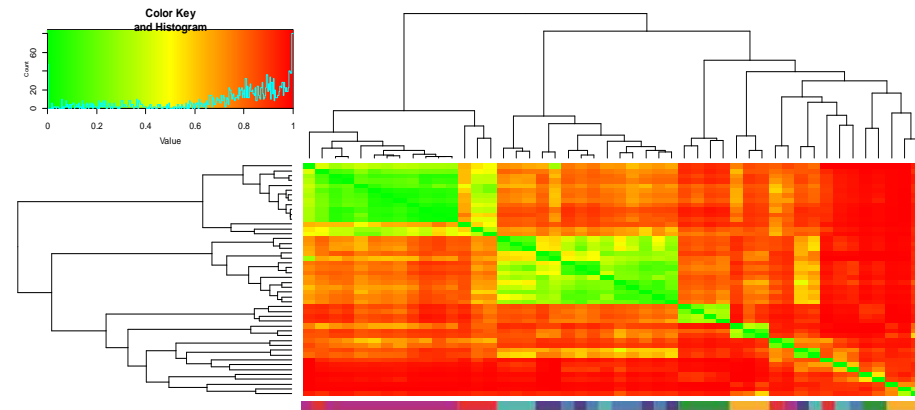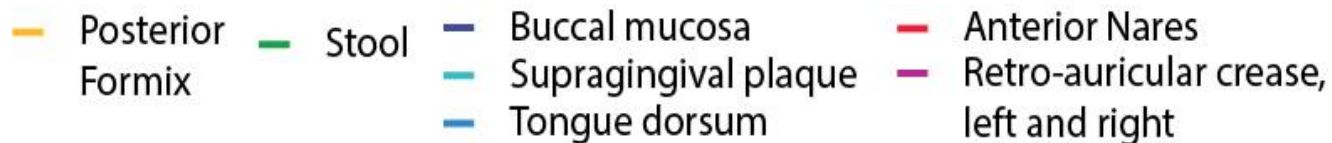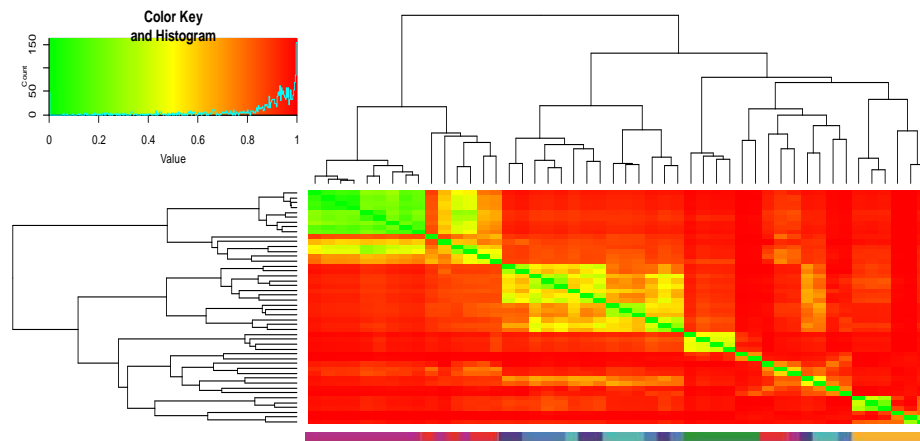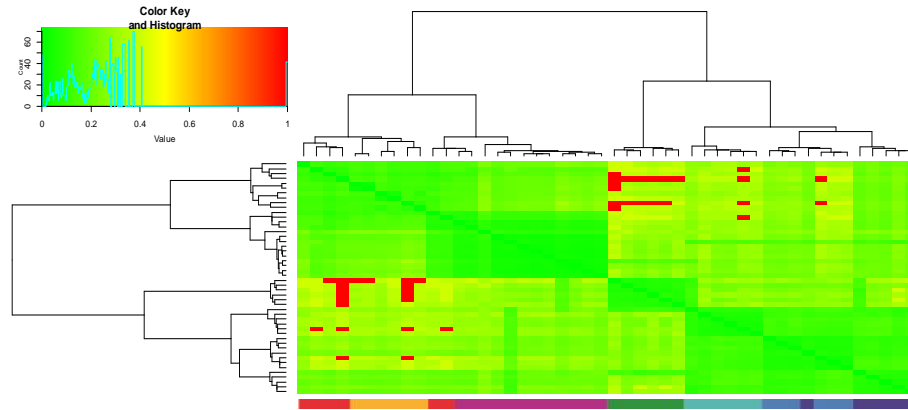
Figure 4
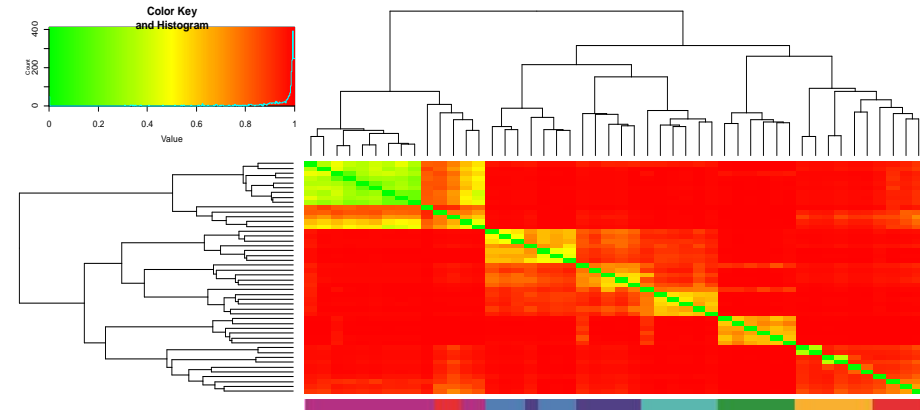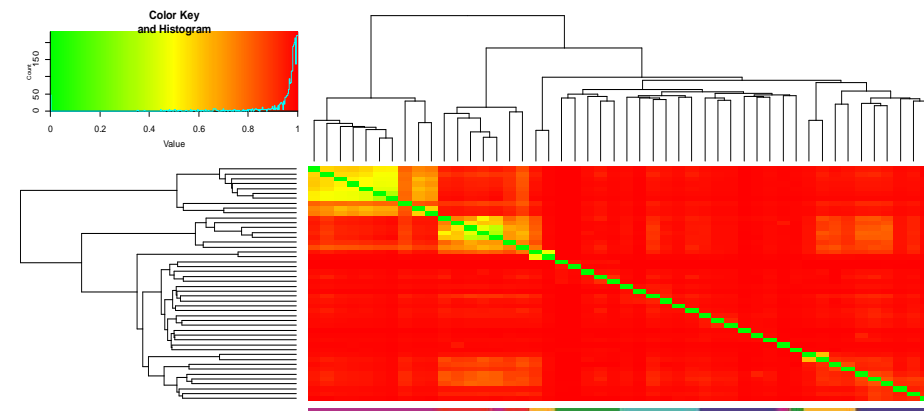
a - MASH

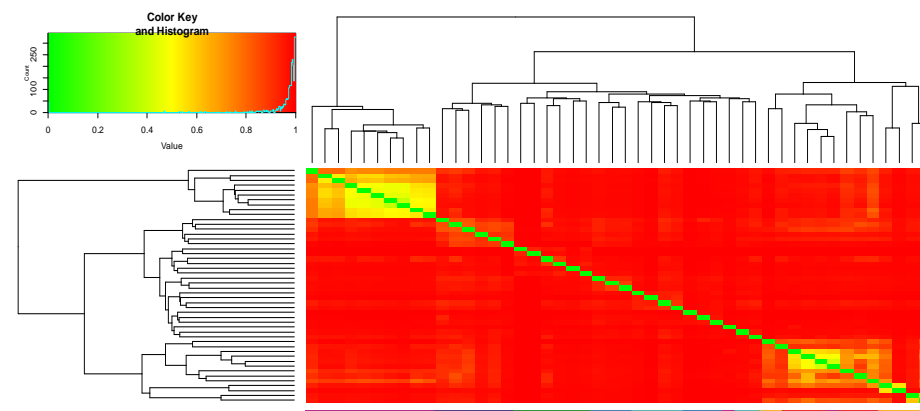b – LIBRA, log weighting

c- SIMKA, abundance Jaccard

d- SIMKA, abundance Bray-Curtis

Posterior Formix

Stool

Buccal mucosa

Supragingival plaque

Tongue dorsum

Anterior Nares

Retro-auricular crease, left and right

Figure 5

C0 - 0~14°C
C1 - 15~21°C
W0 - 22~25°C
W1 - 26~30°C

SUR - surface
DCM - deep chlorophyll maximum
MES - mesopelagic

*edges < 30% are cut



| Mediterranean Sea | South Pacific Ocean |
| Red Sea | North Pacific Ocean |
| Indian Ocean | Southern Ocean |
| South Atlantic Ocean | |

Similarity

0%  10%  20%  30%  40%  50%  60%  70%  80%  90%  100%

Runtimes of Libra

Runtimes of Libra

% to total runtime

6.85%    9.30%    12.64%

0:51:37    1:11:54    1:41:32

11:41:27    11:41:27    11:41:27

Runtime (in hours)

Cosine Similarity    Bray-Curtis    Jensen-Shannon

Distance metrics

Figure

■ index construction    ■ distance-matrix computation

Click here to access/download

**Supplementary Material**

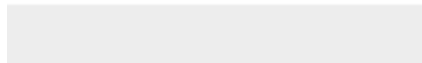Supplemental_methods_and_fig_table_legends.docx

Click here to access/download
**Supplementary Material**
supplemental Figure 1.pdf

Click here to access/download
**Supplementary Material**
Supplemental Figure 2.pdf

Click here to access/download
**Supplementary Material**
Supplemental Figure 3.pdf

Supplementary Table 1

Click here to access/download
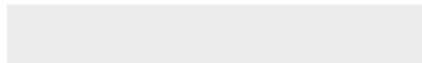**Supplementary Material**
Supplemental Table 1.xlsx

Click here to access/download
**Supplementary Material**
Supplemental Table 2.xlsx

Click here to access/download
**Supplementary Material**
Supplemental Table 3.xlsx

Click here to access/download
**Supplementary Material**
Supplemental Table 4.xlsx

THE UNIVERSITY
OF ARIZONA®

College of Agriculture
and Life Sciences

Department of
Agricultural and
Biosystems
Engineering

Shantz Bldg., B38, Room 403
1177 E. 4th Street
P.O. Box 210038
Tucson, AZ  85721-0038
Tel: (520) 621-1607
Fax: (520) 621-3963

August 24, 2018

Dear Editors,

Please find our paper for consideration at *Gigascience* as a research article titled "Libra: robust biological inferences of global datasets using scalable k-mer based all-vs-all metagenomic comparisons".

Microbiome research spans a broad array of disciplines from medicine, agriculture, bioenergy, and the environment, and is united in addressing core scientific questions relating microbial communities to biological and chemical processes in human, animal, or Earth systems. Given the preponderance of genomic data from diverse environments, there is a new desire to ask cross-cutting questions from the environment to human health. To move this work forward, microbiome datasets need to be holistically analyzed to examine how microbes move through living systems. Currently, only a subset of tools are available that make these analyses possible (through data reduction techniques and read count normalization), but none exploit big data architectures to scale compute and analyze complete datasets (100% of reads) in a linear and fault tolerant manner. This level of resolution is vital in metagenomic analyses where > 50% of the reads are unknown and the only way to understand functional changes in microbial communities is through all-vs-all analysis of diverse datasets to associate sequence patterns with environmental factors. To date, no tool offers a scalable and complete analysis of reads to explore global patterns in microbiome sciences.

Here we describe the first scalable algorithm for comparative metagenomics called Libra that is capable of performing an all-vs-all sequence analysis on hundreds of metagenomes in a Hadoop big data framework. Libra performs with unparalleled accuracy compared to equivalent tools using both simulated and real metagenomic datasets ranging from 80 million to 4.2 billion reads. In contrast to current methods, Libra's state-of-the-art algorithm and its implementation in a big data architecture does not require a reduction in dataset size or simplified distance metrics to achieve remarkable compute times and accuracy. As a result, Libra enables integration of massive datasets across disciplines to identify microbial and viral signatures linked to key biological processes. Moreover, Libra is available as an open-access web-based tool in iMicrobe (http://imicrobe.us) and in Github where the code is available for further optimization and reuse by the community. All authors declare no competing interests and have approved the manuscript for submission. The content of the manuscript has not been published, or submitted for publication elsewhere. Thank you for considering our paper for publication in *Gigascience*.

Sincerely,

Bonnie Hurwitz, PhD
Assistant Professor of Biosystems Engineering
University of Arizona, bhurwitz@email.arizona.edu