

GigaScience

Libra: scalable k-mer based tool for massive all-vs-all metagenome comparisons

--Manuscript Draft--

Manuscript Number:	GIGA-D-18-00324R2	
Full Title:	Libra: scalable k-mer based tool for massive all-vs-all metagenome comparisons	
Article Type:	Technical Note	
Funding Information:	Directorate for Computer and Information Science and Engineering (1640775)	Prof. Bonnie L Hurwitz
Abstract:	<p>Background: Shotgun metagenomics provides powerful insights into microbial community biodiversity and function. Yet, 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.</p> <p>Results: We developed a tool called Libra that performs all-vs-all comparison of metagenomes for precise clustering based on their k-mer content. Libra uses a scalable Hadoop framework for massive metagenome comparisons, Cosine Similarity for calculating the distance using sequence composition and abundance while normalizing for sequencing depth, and a web-based implementation in iMicrobe (http://imicrobe.us) that uses the CyVerse advanced cyberinfrastructure to promote broad use of the tool by the scientific community.</p> <p>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 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 resolution of large-scale comparative analyses. In contrast, Libra uses all of the reads to calculate k-mer abundance in a Hadoop architecture that can scale to any size dataset to enable global-scale analyses and link microbial signatures to biological processes.</p>	
Corresponding Author:	Bonnie Hurwitz 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, BS Ken Youens-Clark, BA John H. Hartman, PhD Bonnie L Hurwitz, PhD	
Order of Authors Secondary Information:		
Response to Reviewers:	GIGA-D-18-00324R1 Libra: scalable k-mer based tool for massive all-vs-all metagenome comparisons Illyoung Choi, MS; Alise J. Ponsero, PhD; Matthew Bomhoff, BS; Ken Youens-Clark, BA; John H. Hartman, PhD; Bonnie L Hurwitz, PhD GigaScience	

Dear Prof. Hurwitz,

Your manuscript "Libra: scalable k-mer based tool for massive all-vs-all metagenome comparisons" (GIGA-D-18-00324R1) has been assessed by our reviewers. Based on these reports, and my own assessment as Editor, I am pleased to inform you that it is potentially acceptable for publication in GigaScience, once you have carried out some essential revisions suggested by our reviewers.

Reviewer #1 requires a few more clarifications to be made. Their reports, together with any other comments, are below. Please also take a moment to check our website at <https://giga.editorialmanager.com/> for any additional comments that were saved as attachments.

RESPONSE: We thank the reviewers for these important additional comments and clarifications. We agree with the reviewers and have addressed the comments in the manuscript per their recommendations. A point-by-point response is provided below. We also ask the the editor consider our resubmission for an Application Note and not a Research Article per reviewer 1's comments below.

In addition, please register any new software application in the SciCrunch.org database to receive a RRID (Research Resource Identification Initiative ID) number, and include this in your manuscript. This will facilitate tracking, reproducibility and re-use of your tool.

RESPONSE: Thank you. We have registered Libra as a tool in SciCrunch.org and have added the RRID (SCR_016608) to the manuscript for tracking and re-use of our tool.

The due date for submitting the revised version of your article is 06 Feb 2019.

We look forward to receiving your revised manuscript soon.

Best wishes,

Nicole Nogoy, Ph.D
GigaScience
www.gigasciencejournal.com

Reviewer reports:

Reviewer #1: Repeating my original observations, Libra appears to be useful and well architected. An extensive comparison to other tools is presented. I appreciate that the authors made specific revisions to the text. However, I feel my most important suggestions were not addressed.

My main suggestion was that this would be better presented as an Application Note, possibly in a different journal. In their response to reviewers, and in defense of submitting a GigaScience Research Article, the authors pointed to their finding that viral communities in the Tara ocean data are similar across temperature gradients, saying this fact was missed in the earlier Tara publication and is being reported here for the first time. If this were the critical finding, then I'd expect it to appear prominently. In fact, it is mentioned twice. First, "Taken together, these data indicate that viral populations are structured globally by temperature, and at finer resolution by station (for surface and DCM samples) indicating that micronutrients and local conditions play an important role in defining viral populations." Second, "We show for the first time that viral communities in the ocean are similar across temperature gradients, irrespective of their location in the ocean."

This treatment does not point out any contradiction to the previous study. The finding is not mentioned in the heading of the subsection, the caption of Table 1 about Tara run time, or the caption of Figure 5 about Tara results. The finding is not mentioned in the Title or in the Abstract or in the Innovations section. The finding appears to be based on a visual interpretation that is vague ("largely structured by temperature") and

provided without statistics. Thus, the wording of the manuscript suggests that this finding was presented, not as a conclusion about the oceans, but as an example of how Libra can be used. In its guide for authors, GigaScience says, "Research Articles present work utilising large scale data that provide some scientific insight and conclusions" (<https://academic.oup.com/gigascience/pages/research>). With respect, I maintain that the revised manuscript is an Application Note and not a Research Article.

RESPONSE: We thank the reviewer for their comments, and agree that the scientific findings are not the main focus of the paper. We ask that the editor consider our revision for an Application Note and not a Research Article.

Secondly, I had noted that the manuscript makes 3 claims to innovation with insufficient support. In their response to reviewers, the authors added the qualification that their application of Hadoop was a first in metagenomics. However, the revised manuscript omits that qualification. After saying, "Libra presents three main innovations", the revised text claims (1) "the use of a scalable Hadoop framework enabling massive dataset comparison" is novel. This sentence does not include any first-in-metagenomics qualification. The claim is unsupported as written.

The revised text claims (2) "linear calculations for complex distance metrics allowing for high accuracy and clustering of the metagenomes based on their k-mer content" is novel. This sentence combines 6 ideas, leaving it unclear what precisely is being claimed. Is this the first linear-time calculation, or the first highly-accurate calculation, or the first k-mer based calculation, or some combination? I find this claim unsupportable as written. The revised text claims (3) "a web-based tool imbedded in the CyVerse advanced cyberinfrastructure through iMicrobe for broader use of the tool in the scientific community" is novel. This claim has no first-in-metagenomics qualification. The claim is unsupported as written. With respect, I maintain that the revised manuscript's three claims to innovation are unproven.

RESPONSE: We agree with the reviewer that each of these claims requires clarification and support based on previous work. The innovation we are trying to convey is in the end-to-end solution we provide rather than each component individually. We have carefully re-phased the abstract and "Innovations" section to clarify this important point. We also added more references and contrasts to previous related works.

We changed the problematic first claim from "the use of a scalable Hadoop framework enabling massive dataset comparison" to "Libra is therefore the first k-mer based de-novo comparative metagenomic tool that uses rely on a Hadoop framework for scalability and fault tolerance"

We changed the second claim from "linear calculations for complex distance metrics allowing for high accuracy and clustering of the metagenomes based on their k-mer content" to "Cosine similarity, although extensively used in computer science, has been rarely implemented in genomic and metagenomic studies (Okuda et al. 2012). To our knowledge, this work is the first to describe the use of the cosine similarity metric to cluster metagenomes based on their k-mer content. "

We modified the last claim from "a web-based tool imbedded in the CyVerse advanced cyberinfrastructure through iMicrobe for broader use of the tool in the scientific community" to "The work described here is the first step in implementing a free cloud-based computing resource for de-novo comparative metagenomics that can be broadly used by scientists to analyze large-scale shared data resources."

A more thorough review might have been possible had Tracked Changes been presented.

RESPONSE: We apologize for the oversight. We have included the tracked changes in three supplemental documents. The first two were from the original re-submission. And the third revision highlights changes described here.

Reviewer #3: Authors have partially addressed my concerns, otherwise several still apply:

The reference 4 that authors give for "Microbial dark matter" does not introduce

	<p>anything about microbial dark matter. Typo ?</p> <p>RESPONSE: Thank you for catching this, we have updated to add three references specific to microbial dark matter and the role of metagenomics in expanding the tree of life.</p> <p>Also, note that it was not necessary to move table 1 to supplemental material -- I was hoping for some clarifications about it not more (cf. my previous comment) -- if authors move this table then they will make sure credits/citations are nonetheless fully given.</p> <p>RESPONSE: To streamline the introduction, we followed your initial suggestion to add the table to the supplemental. We have split the original table into two tables that are focused on the main points in the introduction. Supplemental Table 1A provides a comprehensive list of all de novo metagenomic comparison tools that we are aware of. Supplemental Table 1B provides a comprehensive list of all genomic/metagenomic tools that use a Hadoop framework for computation. The main point of Supplemental Table 1A is to show that Libra is the first de novo metagenomic comparison tool to use a Hadoop framework and also provide the user with a web-based tool. The main point of Supplemental Table 1B is to show that other genomic and metagenomic tools use Hadoop framework, but are for other use-cases. We have also made sure that each of the tools are cited in the main text.</p> <p>There is still the issue of the formatting for the equations/formulas/vectors, see "Cosine Similarity metric" or "Sweep line algorithm", some strange symbols are indicated (I opened this manuscript with different PDF readers, including Adobe, they all show formatting issues). Is it an issue by the editor/s platform or authors ?</p> <p>RESPONSE: Our apologies the conversion didn't work properly again. We fixed by uploading the PDF of the manuscript (as a primary file), in addition to the docx (as Supplemental).</p> <p>Finally, "artificial" is still use in Supplemental Figure 1.</p> <p>RESPONSE: Thank you for finding this. We have updated Supplemental Figure 1 to remove the term "artificial".</p>
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>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.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
Resources	Yes

<p>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.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	
<p>Availability of data and materials</p> <p>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 a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>



[Click here to view linked References](#)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Title: Libra: scalable k-mer based tool for massive all-vs-all metagenome comparisons.

Authors: Illyoung Choi¹, Alise J. Ponsero², Matthew Bomhoff², Ken Youens-Clark², John H. Hartman^{1*}, and Bonnie L. Hurwitz^{2,3*}

ORCID:

Illyoung Choi (<https://orcid.org/0000-0002-9705-6355>)

Alise J. Ponsero (<https://orcid.org/0000-0002-4125-7561>)

Matthew Bomhoff (<https://orcid.org/0000-0002-8014-9184>)

Ken Youens-Clark (<https://orcid.org/0000-0001-9961-144X>)

John H. Hartman (<https://orcid.org/0000-0002-5557-0604>)

Bonnie L. Hurwitz (<https://orcid.org/0000-0001-8699-957X>)

Affiliations:

¹Department of Computer Science, University of Arizona, Tucson, Arizona

²Department of Biosystems Engineering, University of Arizona, Tucson, Arizona

³BIO5 Institute, University of Arizona, Tucson, Arizona

Corresponding Author:

Bonnie L. Hurwitz bhurwitz@email.arizona.edu

1
2
3
4 **ABSTRACT**
5
6

7 **Background:** Shotgun metagenomics provides powerful insights into microbial community
8 biodiversity and function. Yet, inferences from metagenomic studies are often limited by dataset
9 size and complexity and are restricted by the availability and completeness of existing
10 databases. *De novo* comparative metagenomics enables the comparison of metagenomes
11 based on their total genetic content.
12
13
14
15
16
17

18
19
20 **Results:** We developed a tool called Libra that performs an all-vs-all comparison of
21 metagenomes for precise clustering based on their k-mer content. Libra uses a scalable
22 Hadoop framework for massive metagenome comparisons, Cosine Similarity for calculating the
23 distance using sequence composition and abundance while normalizing for sequencing depth,
24 and a web-based implementation in iMicrobe (<http://imicrobe.us>) that uses the CyVerse
25 advanced cyberinfrastructure to promote broad use of the tool by the scientific community.
26
27
28
29
30
31
32

33
34
35 **Conclusions:** A comparison of Libra to equivalent tools using both simulated and real
36 metagenomic datasets, ranging from 80 million to 4.2 billion reads, reveals that methods
37 commonly implemented to reduce compute time for large datasets—such as data reduction,
38 read count normalization, and presence/absence distance metrics—greatly diminish the
39 resolution of large-scale comparative analyses. In contrast, Libra uses all of the reads to
40 calculate k-mer abundance in a Hadoop architecture that can scale to any size dataset to
41 enable global-scale analyses and link microbial signatures to biological processes.
42
43
44
45
46
47
48
49
50

51
52 **Keywords:** metagenomics, Hadoop, k-mer, distance metrics, clustering
53
54
55
56
57
58
59
60
61
62
63
64
65

INTRODUCTION

Over the last decade, scientists have generated petabytes of genomic data to uncover the role of microbes in dynamic living systems. Yet to understand the underlying biological principles that guide the distribution of microbial communities, massive ‘omics datasets need to be compared with environmental factors to find linkages across space and time. One of the greatest challenges in these endeavors has been in documenting and analyzing unexplored genetic diversity in wild microbial communities. For example, fewer than 60% of 40 million non-redundant genes from the Global Ocean Survey (GOS) and the Tara Oceans Expeditions match known proteins in bacteria [1,2]. Other microorganisms such as viruses or pico-eukaryotes that are important to ocean ecosystems are even less well defined (e.g. < 7% of reads from viromes match known proteins [3]). This is largely due to the fact that these organisms are unculturable and reference genomes do not exist in public data repositories. Thus, genome-sequences from metagenomic data await better taxonomic and functional definition. Consequently, even advanced tools such as k-mer based classifiers that rapidly assign metagenomic reads to known microbes miss “microbial dark matter” that comprises a significant proportion of metagenomes [4–6].

De novo comparative metagenomics offers a path forward. In order to examine the complete genomic content, metagenomic samples can be compared using their sequence signature (or frequency of k-mers) (list of tools available in Supplemental Table 1A). This approach relies on three core tenets of k-mer-based analytics: (i) closely related organisms share k-mer profiles and cluster together, making taxonomic assignment unnecessary [7,8], (ii) k-mer frequency is correlated with the abundance of an organism [9], and (iii) k-mers of sufficient length can be used to distinguish specific organisms [10]. In 2012, Compareads [11]

1
2
3
4 method was proposed, followed by Commet [12]. Both of these tools compute the number of
5
6 shared reads between metagenomes using a k-mer-based read similarity measure. The number
7
8 of shared reads between datasets is then used to compute a Jaccard distance between
9
10 samples.

11
12
13 Given the computational intensity of all-vs-all sequence analysis, several other methods have
14
15 been employed to reduce the dimensionality of metagenomes and speed up analyses by
16
17 creating unique k-mer sets and computing the genetic distance between pairs of metagenomes,
18
19 such as MetaFast [13] and Mash [14]. The fastest of these methods, Mash [15], indexes
20
21 samples by unique k-mers to create size-reduced sketches, and compares these sketches
22
23 using the MinHash algorithm [16] for computing a genetic distance using Jaccard similarity. Yet,
24
25 the tradeoff for speed is that samples are reduced to a subset of unique k-mers (1k by default)
26
27 that may lead an unrepresentative k-mer profile of the samples. Further, given that Mash uses
28
29 Jaccard similarity only the genetic distance between samples is accounted for (or genetic
30
31 content in microbial communities) without considering abundance (dominant vs rare organisms
32
33 in the sample) which is central to microbial ecology and ecosystem processes [17]. Sourmash
34
35 [18], a toolkit for manipulating MinHash sketches, uses the same underlying algorithm and
36
37 distance metric as Mash and therefore has the same limitations.
38
39
40
41
42

43 Recently, Simka[15] was developed to compute a distance matrix between metagenomes by
44
45 dividing the input datasets into abundance vectors from subsets of k-mers, then rejoining the
46
47 resulting abundances in a cumulative distance matrix. The methodology can be parallelized to
48
49 execute the analyses on a high-performance computing cluster (HPC). Simka also provides
50
51 various ecological distance metrics to let the user choose the metric most relevant to their
52
53 analysis. However, the computational time varies based on the distance metric, where some
54
55 distances scale linearly and other distances metrics, like Jensen-Shannon, scale quadratically
56
57
58
59
60
61
62
63
64
65

1
2
3
4 as additional samples are added [15]. Moreover, Simka normalizes datasets in an all-vs-all
5
6 comparison by reducing the depth of sequencing for all samples to the least common
7
8 denominator, therefore decreasing the resolution of the datasets. Lastly, computing k-mer
9
10 analytics using HPC is subject to reduced fault tolerance for massive datasets. A framework to
11
12 compare one metagenome to a set of metagenomes on a high-performance computing system
13
14 called DSM [19] has also been proposed, however, this tool is limited to retrieval tasks and does
15
16 not provide an all-vs-all sequence analysis.
17
18

19
20 **Scaling sequence analysis using big data analytics via Hadoop.** Hadoop is an attractive
21
22 platform for performing large-scale sequence analysis because it provides a distributed file
23
24 system and distributed computation for analyzing massive amounts of data. Hadoop clusters are
25
26 comprised of commodity servers so that the processing power increases as more computing
27
28 resources are added. Hadoop also offers a high-level programming abstraction, called
29
30 MapReduce [20] that greatly simplifies the implementation of new analytical tools and a
31
32 high-performance distributed file system (HDFS) for storing data sets. Programmers do not
33
34 need specialized training in distributed systems and networking to implement distributed
35
36 programs using MapReduce. Hadoop also provides fault-tolerance by default. When a Hadoop
37
38 node fails, Hadoop reassigns the failed node's tasks to another node containing a redundant
39
40 copy of the data those tasks were processing. This differs from HPC where schedulers track
41
42 failed nodes and either restart the failed computation from the most recent checkpoint, or from
43
44 the beginning if checkpointing wasn't used. Thus, using a Hadoop infrastructure ensures that
45
46 computations and data are protected even in the event of hardware failures. These benefits
47
48 have led to new analytic tools based on Hadoop, making Hadoop a de facto standard in
49
50 large-scale data analysis. In metagenomics, the development of efficient and inexpensive
51
52 high-throughput sequencing technologies has lead to a rapid increase in the amount of
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 sequence data for studying microbes in diverse environments. However, to date only
5
6 Hadoop-enabled genomic or k-mer counting tools exist, and no comparative metagenomics
7
8 tools are available (Supplemental Table 1B).
9

10
11 **Existing big data algorithms compare reads to limited genomic reference data.** Recent
12
13 progress has been made in translating bioinformatics algorithms to big data architectures to
14
15 overcome scalability issues. Thus far, these algorithms compare large-scale NGS datasets to
16
17 reference genomic datasets and replace computationally intensive algorithms such as sequence
18
19 alignment [21], genetic variant detection [22,23], ortholog detection[24], differential gene
20
21 expression[25,26], or short read mapping [27–30] (Supplemental Table 1B). For example,
22
23 BlastReduce and CloudBurst are parallel sequence mapping tools based on Hadoop
24
25 MapReduce [28,29]. These tools, however, implement a query-to-a-reference approach that is
26
27 inefficient for all-vs-all analyses of reads from metagenomes. Other algorithms such as BioPig
28
29 [31] and Bloomfish [32] generate an index of sequence data for later partial sequence search
30
31 and k-mer counting using Hadoop [33] (Supplemental Table 1B). Also, some of these tools
32
33 adopt traditional sequence indexing techniques such as a suffix array that is inefficient in
34
35 reading and indexing data in HDFS, thus reducing performance. Moreover, neither tool offers an
36
37 end-to-end solution for comparing metagenomes consisting of data distribution on a Hadoop
38
39 cluster, k-mer indexing and counting, distance matrix computation, and visualization. Finally,
40
41 none of these tools are enabled in an advanced cyberinfrastructure where users can compute
42
43 analyses in a simple web-based platform (Supplemental Table 1B).
44
45
46
47
48
49

50
51 **Libra: a tool for scalable all-vs-all sequence analysis in an advanced cyberinfrastructure**

52
53 Here, we describe a scalable algorithm called Libra that is capable of performing all-vs-all
54
55 sequence analysis using Hadoop MapReduce (SciCrunch.org tool reference ID SCR_016608).
56
57 We demonstrate for the first time that Hadoop MapReduce can be applied to all-vs-all sequence
58
59
60
61
62
63
64
65

1
2
3
4 comparisons of large-scale metagenomic datasets comprised of mixed microbial communities.
5
6 We demonstrate that Cosine Similarity, which is widely used in document clustering and
7
8 information retrieval, is a good distance metric for comparing datasets to consider genetic
9
10 distance and microbial abundance simultaneously, along with widely accepted distance metrics
11
12 in biology such as Bray-Curtis [34] and Jensen-Shannon [35]. We validate this distance metric
13
14 using simulated metagenomes (from both short and long read technologies) to show that Libra
15
16 has exceptional sensitivity in distinguishing complex mixed microbiomes. Next, we show Libra's
17
18 ability to distinguish metagenomes by both community composition and abundance using 48
19
20 samples (16S rRNA and WGS) from the human microbiome project (HMP) and the simulated
21
22 Critical Assessment of Metagenome Interpretation (CAMI) "toy" PacBio dataset across diverse
23
24 body sites and compare the results to Mash and Simka. Finally, we show that Libra can scale to
25
26 massive global-scale datasets by examining viral diversity in 43 Tara Ocean Viromes (TOV)
27
28 from the 2009-2011 Expedition [36] that represent 26 sites containing about 4.2 billion reads.
29
30 We show for the first time that viral communities in the ocean are similar across temperature
31
32 gradients, irrespective of their location in the ocean. The resulting data demonstrate that Libra
33
34 provides accurate, efficient, and scalable computation for comparative metagenomics that can
35
36 be used to discern global patterns in microbial ecology.
37
38
39
40
41
42
43

44 To promote the broad use of the Libra algorithm we developed a web-based tool in iMicrobe
45
46 [37], where users can run Libra using data in their free CyVerse [38,39] account or use datasets
47
48 that are integrated into the iMicrobe Data Commons. These analyses are fundamental for
49
50 determining relationships among diverse metagenomes to inform follow-up analyses on
51
52 microbial-driven biological processes.
53
54
55
56

57 **DATA DESCRIPTION**

58
59
60
61
62
63
64
65

1
2
3
4 **Staggered mock community.** We performed metagenomic shotgun sequencing on a
5
6 staggered mock community obtained from the Human Microbiome Consortium (HM-277D). The
7
8 staggered mock community is comprised of genomic DNA from genera commonly found on or
9
10 within the human body, consisting of 1,000 to 1,000,000,000 16S rRNA gene copies per
11
12 organism per aliquot. The resulting DNA was subjected to whole genome sequencing as
13
14 follows. Mixtures were diluted to a final concentration of 1 nanogram/microliter and used to
15
16 generate whole genome sequencing libraries with the Ion Xpress Plug Fragment Library Kit and
17
18 manual #MAN0009847, revC (Thermo Fisher Scientific, Waltham, MA, USA). Briefly, 10
19
20 nanograms of bacterial DNA was sheared using the Ion Shear enzymatic reaction for 12 min
21
22 and Ion Xpress barcode adapters ligated following end repair. Following barcode ligation,
23
24 libraries were amplified using the manufacturer's supplied Library Amplification primers and
25
26 recommended conditions. Amplified libraries were size selected to ~ 200 base pairs using the
27
28 Invitrogen E-gel Size Select Agarose cassettes as outlined in the Ion Xpress manual and
29
30 quantitated with the Ion Universal Library quantitation kit. Equimolar amounts of the library were
31
32 added to an Ion PI Template OT2 200 kit V3. The resulting templated beads were enriched with
33
34 the Ion OneTouch ES system and quantitated with the Qubit Ion Sphere Quality Control kit (Life
35
36 Technologies) on a Qubit 3.0 fluorometer (Qubit, NY, NY, USA). Enriched templated beads
37
38 were loaded onto an Ion PI V2 chip and sequenced according to the manufacturer's protocol
39
40 using the Ion PI Sequencing 200 kit V3 on an Ion Torrent Proton sequencer. The sequence data
41
42 comprised of ~80 million reads have been deposited to the NCBI Sequence Read Archive under
43
44 accession SRP115095 under project accession PRJNA397434.
45
46
47
48
49
50
51

52 **Simulated data derived from the staggered mock community.** The resulting sequence data
53
54 from the staggered mock community (~80 million reads) were used to develop simulated
55
56 metagenomes to test the effects of varying read depth, and composition and abundance of
57
58
59
60
61
62
63
64
65

1
2
3
4 organisms in mixed metagenomes [40]. To examine read depth (in terms of raw read counts
5 and file size), we used the known staggered mock community abundance profile to generate a
6 simulated metagenome using GemSim [41] of 2 million reads (454 sequencing) and duplicated
7 the dataset 2x, 5x and 10x. We also simulated the effects of sequencing a metagenome more
8 deeply using GemSim [41] to generate simulated metagenomes with 0.5, 1, 5, and 10 million
9 reads based on the relative abundance of organisms in the staggered mock community. Next,
10 we developed four simulated metagenomes to test the effect of changing the dominant
11 organism abundance and genetic composition including: 10 million reads from the staggered
12 mock community (mock 1), the mock community with alterations in a few abundant species
13 (mock 2), the mock community with many alterations in abundant species (mock 3), and mock 3
14 with additional sequences from archaea to further alter the genetic composition (mock 4) as
15 described in Supplemental Table 2. The same community profiles were used to generate
16 paired-end Illumina dataset (100 million reads), using GemSim (Illumina v4 error model). Finally,
17 using SimLord [42], the community profiles were used to generate simulated third-generation
18 sequencing datasets (Pacific Bioscience SMRT sequencing - 1 million reads). SimLord default
19 parameters were used to generate those simulated datasets. All simulated datasets are
20 available in iMicrobe [37] under project 265 and under DOI[40].
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43

44 **Human microbiome 16S rRNA gene amplicons and WGS reads.** Human microbiome
45 datasets were downloaded from the NIH Human Microbiome Project [43] including 48 samples
46 from 5 body sites including: urogenital (posterior fornix), gastrointestinal (stool), oral (buccal
47 mucosa, supragingival plaque, tongue dorsum), airways (anterior nares), and skin
48 (retroauricular crease left and right) ([See Supplemental Table 3]). Matched datasets consisting
49 of 16S rRNA reads WGS reads, and WGS assembled contigs were downloaded from the 16S
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 trimmed dataset and the HMIWGS/HMASM dataset respectively. For the WGS reads dataset,
5
6 the analysis was run on the paired 1 read file.
7
8

9
10 **Tara ocean viromes.** Tara oceans viromes were downloaded from European Nucleotide
11
12 Archive (ENA) at EMBL and consisted of 43 viromes from 43 samples at 26 locations across the
13
14 world's oceans collected during the Tara Oceans (2009-2012) scientific expedition
15
16 (Supplemental Table 4) [36]. Metadata for the samples were downloaded from PANGAEA [44].
17
18 These samples were derived from multiple depths including 16 surface samples (5-6 meters),
19
20 18 deep chlorophyll maximum samples (DCM; 17-148 meters), and one mesopelagic sample
21
22 (791 meters). Quality control procedures were applied according to the methods described by
23
24 Brum and colleagues [36].
25
26
27

28 29 **CAMI Human microbiome project toy dataset**

30
31
32 The human microbiome project toy dataset from the Critical Assessment of Metagenome
33
34 Interpretation (CAMI) 2nd Challenge was downloaded from their website [45]. This dataset is
35
36 composed of 49 simulated PacBio reads from five different body sites of the human host,
37
38 namely gastrointestinal tract, oral cavity, airways, skin and urogenital tract.
39
40
41

42 **RESULTS AND DISCUSSION**

43
44
45 **Libra computational strategy.** Libra uses Hadoop MapReduce to perform massive all-vs-all
46
47 sequence comparisons between next-generation sequence (NGS) datasets. Libra uses a
48
49 scalable algorithm and efficient resource usage to make all-vs-all comparisons feasible on large
50
51 datasets. Hadoop allows parallel computation over distributed computing resources via its
52
53 simple programming interface called *MapReduce*, while hiding much of the complexity of
54
55 distributed computing (e.g. node failures) for robust fault-tolerant computation. Taking
56
57
58
59
60
61
62
63
64
65

1
2
3
4 advantage of Hadoop, Libra can scale to larger input datasets and more computing resources.
5
6 Furthermore, many cloud providers such as Amazon and Google offer Hadoop clusters on a
7
8 pay-as-you-go basis, allowing scientists to scale their Libra computations to match their
9
10 datasets and budgets.
11
12

13
14 Libra is implemented using three different MapReduce jobs — 1) k-mer histogram construction,
15
16 2) inverted index construction, and 3) distance matrix computation. Fig 1 shows a workflow of
17
18 the Libra algorithm.
19
20
21

22 **Figure 1. The Libra Workflow.**

23
24

25
26 Libra consists of three MapReduce jobs (yellow boxes) — 1) Libra constructs a k-mer histogram
27
28 of the input samples for load-balancing. The k-mer histogram of the input samples is computed
29
30 in parallel by running multiple Map tasks and a Reduce task that combines their results; 2) Libra
31
32 constructs the inverted index in parallel. In the Map phase, a separate Map task is spawned for
33
34 every data block in the input sample files. Each Map task generates k-mers from the sequences
35
36 stored in a data block then passes them to the Reduce tasks. Each Reduce task then counts
37
38 k-mers it receives and produces an index chunk; 3) In the distance matrix computation, the work
39
40 is split by partitioning the k-mer space at the beginning of a MapReduce job. The k-mer
41
42 histogram files for input samples are loaded and the k-mer space is partitioned according to the
43
44 k-mer distributions. A separate Map task is spawned for each partition to perform the
45
46 computation in parallel and merged to produce the complete distance matrix.
47
48
49

50
51 **Libra distance computation.** Jaccard and Bray-Curtis distance have been extensively used to
52
53 compare metagenomes based on their sequence signature [13–15]. While Mash only computes
54
55 the Jaccard distance between samples, Simka and Libra implement several classical ecology
56
57 distances, allowing the user to choose the best-suited distance for the considered dataset [15].
58
59
60
61
62
63
64
65

1
2
3
4 Libra provides three distance metrics — Cosine Similarity, Bray-Curtis, and Jensen-Shannon. In
5
6 this paper, we demonstrate Cosine Similarity as the default distance metric. This distance uses
7
8 a vector space model to compute the distance between two NGS samples based on their k-mer
9
10 composition and abundance, while simultaneously normalizing for sequencing depth. Cosine
11
12 Similarity is widely used in document clustering and information retrieval. This distance metric
13
14 was previously used to evaluate the accuracy of methods to reconstruct genomes from “virtual
15
16 metagenomes” derived from 16S rRNA data based on shared KEGG orthologous gene counts
17
18 [46] but has not been applied in analyzing sequence signatures between metagenomes. Libra
19
20 users can also weight k-mers based on their abundance (using boolean weighting, natural
21
22 weighting, and logarithmic weighting) to account for differences in microbial community
23
24 composition and sequencing effort as detailed below.
25
26
27
28
29

30
31 **Cosine Similarity allows for an accurate and normalized comparison of metagenomes.**

32
33 We explored the effects of varying: (1) the size of the datasets, (2) depth of sequencing, (3) the
34
35 abundance of dominant microbes in the community, and (4) genetic composition of the
36
37 community by adding in an entirely new organism (in our case we added archaea). We
38
39 constructed simulated metagenomes and compared Libra’s distance based on the Cosine
40
41 Similarity against those from Mash and Simka. Simulated datasets were derived from genomic
42
43 DNA from a staggered mock community of bacteria obtained from the human microbiome
44
45 consortium and sequenced deeply using the Ion Torrent sequencing platform (80 million reads,
46
47 see methods).
48
49

50
51 First, we examined the effect of the size of the dataset by using GemSim [41] to obtain a
52
53 simulated metagenome composed of 1 million reads (454 sequencing) from the mock
54
55 community and duplicating that dataset 2x and 10x. Overall, we found that altering the size of
56
57 the metagenome (by duplicating the data) had no effect on the distance between metagenomes
58
59
60
61
62
63
64
65

1
2
3
4 for Mash, Simka, or Libra. In each case, the distance of the duplicated datasets to the 1x mock
5
6 community was less than 0.0001 (data not shown).
7
8

9
10 Because metagenomes don't scale exactly with size and instead have an increasing
11
12 representation of low-abundance organisms, we created a second simulated dataset from the
13
14 mock community using GemSim [41] 0.5, 1, 5, and 10 million reads (454 sequencing) to mimic
15
16 the effect of reducing the sequencing. Given the abundance of organisms in the mock
17
18 community, the 0.5 M read dataset is mainly comprised of dominant species. Because Simka
19
20 normalizes all samples to the lowest read count, no changes between samples were
21
22 measurable when using Jaccard and Bray-Curtis distances (Fig 2A). In contrast, Mash and Libra
23
24 (natural weighting) take into account all of the reads in the metagenomes, therefore they
25
26 measure a larger difference when you compare the smallest (0.5M read sample) and largest (10
27
28 million read sample). These results suggest that Libra (natural weighting) and Mash are
29
30 appropriate for comparing datasets at different sequencing depths, whereas using Simka could
31
32 lead to undesired effects.
33
34
35
36
37

38 **Figure 2. Analysis of simulated metagenomes using Mash, Simka, and Libra.**

- 39
40
41
42 A. Distance to staggered mock community simulated metagenome composed of 10 million
43
44 reads (mock1 10M), for simulated metagenomes of same community sequenced at
45
46 various depth. Simulated metagenomes (454 sequencing) were obtained using GemSim
47
48 and the known abundance profile of the staggered mock community (see Supplemental
49
50 Table 2). In order to mimic various sequencing depths, the simulated metagenomes
51
52 were generated at 0.5, 1, 5 or 10 million reads (noted mock1 0.5M; mock1 1M; mock1
53
54 5M; mock1V2 10M). The distances between the 4 simulated metagenomes and a 10
55
56 million read simulated metagenome (mock1 10M) were computed using Mash, Simka
57
58
59
60
61
62
63
64
65

1
2
3
4 (Jaccard and Bray-Curtis distance) and Libra (natural weighting).
5

6
7 B. Distance to staggered mock community simulated metagenome (mock 1), for simulated
8 metagenomes from increasingly distant communities. The mock 1 relies on the known
9 abundance profile from the staggered mock community. The mock 2 community profile
10 was obtained by randomly inverting 3 species abundance from mock 1 profile. The mock
11 3 profile was obtained by randomly inverting 2 species abundances from mock 2 profile.
12
13 Finally, a mock 4 profile was obtained by adding high abundance archeal genomes not
14 present in any the other mock communities. Simulated metagenomes (454 sequencing)
15 were generated using GemSim at 10 million reads. The distance between the mock 1
16 community to mock 2, mock 3, mock 4 and a replicate community (mock1 V2) was
17 computed using Mash, Simka (Jaccard and Bray-Curtis distance) and Libra (cosine
18 distance, natural and logarithmic weighting).
19
20
21
22
23
24
25
26
27
28
29
30
31

32 In addition to natural variation in population-level abundances, artifacts from sequencing can
33 result in high-abundance k-mers. Libra allows users to select the optimal methodology for
34 weighting high abundance k-mers in their datasets including boolean, natural, and logarithmic.
35
36 These options for weighting k-mers are important for different biological scenarios as described
37 below and shown in simulated datasets. To examine the effect of weighting, we compared and
38 contrasted the natural and logarithmic weight in Libra, with other distances obtained from Mash
39 and Simka (Jaccard and Bray-Curtis). We also examined the effect of adding an entirely new
40 species by spiking a simulated dataset with sequences derived from archaea (that were not
41 present in the mock community). The simulated datasets (454 technology) were comprised of
42 the staggered mock community (mock 1), the mock community with alterations in a few
43 abundant species (mock 2), the mock community with many alterations in abundant species
44 (mock 3), and mock 3 with additional sequences from archaea to alter the genetic composition
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 of the community (mock 4) (see Supplemental Table 2). The resulting data showed that Libra
5
6 (logarithmic weighting) shows a stepwise increase in distance among the mock communities
7
8 (Fig 2B). This suggests that logarithmic weighting in Libra allows for a comparison of distantly
9
10 related microbial communities. Mash also shows a stepwise distance between communities but
11
12 is compressed relative to Libra, making differences less distinct. Simka (Bray-Curtis and
13
14 Jaccard) and Libra (cosine distance, natural weighting) reach the maximum difference between
15
16 mock communities 3 and 4 (Fig 2B). This indicates that these distances are more appropriate
17
18 when comparing metagenomes with small fluctuations in the community (e.g., data from a
19
20 time-series analysis), whereas Libra (cosine distance, logarithmic weighting) can be used to
21
22 distinguish metagenomes that vary in both genetic composition and abundance over a wide
23
24 range of species diversity by dampening the effect of high-abundance k-mers. Because of this
25
26 important difference, we used the cosine distance with the logarithmic weighting in all
27
28 subsequent analyses. Further, we also found that cosine distance provides the fastest
29
30 computation among all distance metrics (see Methods). We confirmed these findings using
31
32 Illumina simulated datasets (Supplemental Figure 1A), to show that these results are consistent
33
34 across short-read technologies.
35
36
37
38
39
40
41

42 Given the availability of long read (~10K) sequencing technologies like Oxford Nanopore and
43
44 PacBio sequencing, we repeated the analyses above on simulated long read data
45
46 (Supplemental Figure 1B). We show that simulated PacBio long read data for the mock
47
48 community derived from SimLoRD [42] shows a similar stepwise distance pattern between each
49
50 of the mock communities (Supplemental Figure 1B), but has a higher overall distance between
51
52 mock 1 and each of the mock communities (mock 2 - 4) likely due to the high simulated random
53
54 error rate compared to simulated short read data.
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 **Libra accurately profiles differences in bacterial diversity and abundance in amplicon**
5 **and WGS datasets from the human microbiome.**
6

7
8
9 Microbial diversity is traditionally assessed using two methods: the 16S rRNA gene to classify
10 bacterial and archaeal groups at the genus to species level, or whole genome shotgun
11 sequencing (WGS) for finer taxonomic classification at the species or subspecies level. Further,
12
13 WGS datasets provide additional information on functional differences between metagenomes.
14
15 Here we compare and contrast the effect of different algorithmic approaches (Mash vs Libra vs
16 Simka), distance metric (Libra vs Simka), data type (16S rRNA vs WGS), and sequence type
17
18 (WGS reads vs assembled contigs) in analyzing data from 48 samples across 8 body sites from
19
20 the Human Microbiome Project. Specifically, we examine matched datasets (16S rRNA reads,
21
22 WGS reads, and WGS assembled contigs) classified as urogenital (posterior fornix),
23
24 gastrointestinal (stool), oral (buccal mucosa, supragingival plaque, tongue dorsum), airways
25
26 (anterior nares), and skin (retroauricular crease left and right) ([See Supplemental Table 2]).
27
28
29
30
31

32
33
34
35 Because the HMP datasets represent microbial communities, abundant bacteria will have more
36
37 total read counts than rare bacteria in the samples. Thus, each sample can vary by both taxonomic
38
39 composition (the genetic content of taxa in a sample) and abundance (the relative proportion of
40
41 those taxa in the samples). Importantly, the 16S rRNA amplicon dataset is useful in showing how
42
43 well each algorithm performs in detecting and quantifying small-scale variation for single a gene at
44
45 the genus-level, whereas the WGS dataset demonstrates the effect of including the complete
46
47 genetic content and abundance of organisms at the species-level in a community [47]. Also, we
48
49 examine differences in each algorithm when read abundance is excluded using assembled contigs
50
51 that only represent the genetic composition of the community.
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 Using the 16S rRNA reads, both Mash and Libra clustered samples by broad categories but not
5 individual body-sites (Fig 3A and B). Similar to what is described in previous work [15], samples
6 from the airways and skin co-cluster, whereas other categories including urogenital,
7 gastrointestinal, and oral are distinct [15]. These results indicate that limited variation in the 16S
8 rRNA gene may only allow for clustering for broad categories. Further, the Mash algorithm shows
9 lower overall resolution (Fig 3A) as compared to Libra (Fig 3B). Indeed, amplicon sequencing
10 analysis is not an original intended use of Mash, given that it reduces the dimensionality of the data
11 by looking at presence/absence of unique k-mers, whereas Libra examines the complete dataset
12 accounting for both the genetic composition of organisms and their abundance. In contrast, Simka
13 (Jaccard-ab and Bray-Curtis) fails to cluster samples by broad categories: some skin samples are
14 found associated with stool and fornix samples (Fig 3C and D). Moreover, Simka Jaccard-ab fails
15 to cluster the mouth samples together (Fig 3C). This result suggests that applying Simka and these
16 well-used distance metrics are not appropriate for these datasets.
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

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

36
37 48 Human metagenomic samples from the HMP projects clustered by Mash (A), Libra (B) or
38 Simka using Jaccard-ab (C) and Bray-Curtis distances (D) from 16S rRNA sequencing runs.

39
40 The samples were clustered using Ward's method on their distance scores. Mash, Simka, and
41 Libra report distance in the same range (0-1). Heat maps showing the pairwise dissimilarity
42 between samples were therefore scaled between 0 (green) and 1 (red). A key below the
43 heatmap colors the samples by body sites.
44
45
46
47
48
49
50

51
52 When using WGS reads, both Mash and Libra show enhanced clustering by body-site (Fig 4A and
53 B), however, Mash shows decreased resolution (Fig 4A) as compared to Libra (Fig 4B). Again,
54 these differences reflect the effect of using all of the read data (Libra) rather than a subset (Mash).
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 The effect of using all of the read data compared to a subset (when sketching in Mash) has been
5
6 previously described in Benoit *et al.* [15]. Importantly, the Libra algorithm depends on read
7
8 abundance that provides increased resolution for interpersonal variation as seen in skin samples
9
10 (Fig 4B). Similar to the 16S rRNA datasets, Simka (Jaccard-ab and Bray-Curtis) failed to cluster
11
12 the samples by body site, where some skin and stool samples cluster with fornix samples (Fig 4C
13
14 and D). Similarly, Simka Jaccard-ab also fails to cluster the mouth samples together (Fig 4C).
15
16 Overall Simka shows an enhanced clustering by body-site using WGS data compared to the 16S
17
18 rRNA data using these distance metrics, however, the clustering is still not accurate. In order to
19
20 confirm the independence of these result toward the sequencing technology, we performed the
21
22 same experiment on the *CAMI HMP* “toy dataset” (simulated PacBio long reads) [Supplemental
23
24 Figure 2]. This analysis shows that each of the tools is able to cluster the samples broadly by body
25
26 site. However, there are small misclassifications shared across all tools, suggesting that the
27
28 increased error rate for this technology could have a limited impact on k-mer based analytics.
29
30
31
32
33
34

35 **Figure 4. Clustering of WGS samples using Mash, and Libra and Simka.**

36
37 48 Human metagenomic samples from the HMP projects clustered by Mash (A), Libra (B) or
38
39 Simka using Jaccard-ab (C) and Bray-Curtis distances (D) from whole genome shotgun
40
41 sequencing runs. The samples were clustered using Ward’s method on their distance scores.
42
43 Heat maps illustrate the pairwise dissimilarity between samples, scaled between 0 (green) and
44
45 1 (red). A key below the heatmap colors the samples by body sites.
46
47
48
49

50 When abundance is taken out of the equation by using assembled contigs ([See Supplemental
51
52 Figure 3]) Mash performs well in clustering distinct body sites whereas Libra shows discrepancies
53
54 and less overall resolution. Thus, as designed Libra requires reads rather than contigs to perform
55
56 accurately and obtain high-resolution clustering (Fig 4). Simka (Jaccard-ab and Bray-Curtis) was
57
58
59
60
61
62
63
64
65

1
2
3
4 not able to distinguish any assembled datasets and scored all sample-to-sample distances to the
5
6 maximum, even considering presence-absence distance metric proposed by Simka (data not
7
8 shown). This phenomenon may be explained by the normalization method used by Simka, which
9
10 does not provide enough data to compare the samples when normalized by the smallest number of
11
12 contigs (in our dataset 69 contigs).
13
14

15
16
17 **Libra allows for ecosystem-scale analysis: clustering the Tara ocean viromes to unravel**
18
19 **global patterns.**
20

21 To demonstrate the scale and performance of the Libra algorithm, we analyzed 43 Tara Ocean
22
23 Viromes (TOV) from the 2009-2011 Expedition [36] representing 26 sites, 43 samples, and 4.2
24
25 billion reads from the global ocean (see Methods). Phages (viruses that infect bacteria) are
26
27 abundant in the ocean [48] and can significantly impact environmental processes through host
28
29 mortality, horizontal gene transfer, and host-gene expression. Yet, how phages change over
30
31 space and time in the global ocean and with environmental fluxes is just beginning to be
32
33 explored. The primary challenge is the majority of reads in viromes (often > 90%) do not match
34
35 known proteins or viral genomes [3] and no conserved genes like the bacterial 16S rRNA gene
36
37 exist to differentiate populations. To examine known and unknown viruses simultaneously,
38
39 viromes are best compared using sequence signatures to identify common viral populations.
40
41
42

43
44 Two approaches exist to cluster viromes based on sequence composition. The first approach
45
46 uses protein clustering to examine functional diversity in viromes between sites [3,36,49].
47
48 Protein clustering, however, depends on accurate assembly and gene finding that can be
49
50 problematic in fragmented and genetically diverse viromes [50]. Further, assemblies from
51
52 viromes often include only a fraction of the total reads (e.g., only 1/3 in TOV [36]). To examine
53
54 global viral diversity in the ocean using all of the reads we examined TOV using Libra. The
55
56 complete pairwise analysis of ~4.2 billion reads in the TOV dataset [36] finished in 18 hours
57
58
59
60
61
62
63
64
65

1
2
3
4 using a 10-node Hadoop cluster (see Methods and Table 2). Importantly, Libra exhibits
5
6 remarkable performance in computing the distance matrix, wherein k-mer matches for all TOV
7
8 completed within 1.5 hours (see Table 1). This step usually represents the largest computational
9
10 bottleneck for bioinformatics tools that compute pairwise distances between sequence pairs for
11
12 applications such as hierarchical sequence clustering [51–54]. A direct comparison of the
13
14 runtime of the Simka, Mash, and Libra are not possible given that each tool is tuned to a
15
16 different computational architecture with a different number of servers and total CPU/memory
17
18 (Mash runs on a single server; Simka runs on an HPC, and Libra on Hadoop).
19
20
21
22

23
24 **Table 1. Execution times for the Libra based on the Tara Ocean Virome (TOV) dataset.**
25

27 Stage	28 Execution Time
29 Preprocessing 30 (k-mer histogram construction 31 / Inverted index construction)	32 16:32:55
33 Distance matrix computation	34 1:24:27
35 Total	36 17:57:22

37
38
39
40
41

42 Overall, we found that viral populations in the ocean are largely structured by temperature in
43
44 four gradients (Fig 5) similar to their bacterial hosts [2]. Interestingly, samples from different
45
46 Longhurst Provinces but the same temperature gradient cluster together. Also, water samples
47
48 from the surface (SUR) and deep chlorophyll maximum (DCM) at the same station, cluster more
49
50 closely together than samples from the same depth at nearby sites (Fig 5). Also noteworthy,
51
52 samples that were derived from extremely cold environments (noted as C0 in Fig 5) lacked
53
54 similarity to all other samples (at a 30% similarity score), indicating distinctly different viral
55
56 populations. These samples include a mesotrophic sample that has previously been shown to
57
58
59
60
61
62
63
64
65

1
2
3
4 have distinctly different viral populations than surface ocean samples [55]. Taken together,
5
6 these data indicate that viral populations are structured globally by temperature, and at finer
7
8 resolution by the station (for surface and DCM samples) indicating that micronutrients and local
9
10 conditions play an important role in defining viral populations.
11
12
13

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

15
16 Similarities between samples from 43 TOV from the 2009-2012 Tara Oceans Expedition. Lines
17
18 (edges) between samples represent the similarity and are colored and thickened accordingly.
19
20 Lines with insignificant similarity (less than 30%) are removed. Each of the sample names is
21
22 color-coded by Longhurst Province. Inner circles show temperature ranges. Sample names
23
24 show the temperature range, station, and depth as indicated on the legend. The analysis is
25
26 performed using Libra (k=20, Logarithmic weighting, and Cosine Similarity).
27
28
29
30
31
32

33 **INNOVATIONS**

34
35
36 Scientific collaboration is increasingly data-driven given large-scale next-generation sequencing
37
38 datasets. It is now possible to generate, aggregate, archive, and share datasets that are
39
40 terabytes and even petabytes in size. Scalability of a system is becoming a vital feature that
41
42 decides the feasibility of massive 'omic's analyses. In particular, this is important for
43
44 metagenomics where patterns in global ecology can only be discerned by comparing the
45
46 sequence signatures of microbial communities from massive 'omics datasets, given that most
47
48 microbial genomes have not been defined. Current algorithms to perform these tasks run on
49
50 local workstations or high-performance computing architectures.
51
52
53
54

55
56 Hadoop is a well-used framework allowing for scalability. The Hadoop framework was previously
57
58 used for k-mer spectra calculation in prior work (Supplemental Table 1B) [31][32]. However,
59
60
61
62
63
64
65

1
2
3
4 these tools do not provide any distance computation between the generated k-mer spectra. To
5
6 our knowledge, Libra is, therefore, the first k-mer based *de-novo* comparative metagenomic tool
7
8 that uses a Hadoop framework for scalability and fault tolerance.
9

10
11 *De-novo* comparative metagenomic tools rely on the calculation of a distance metric in order to
12
13 perform a clustering task on the metagenomes. Libra provides several distance metrics on the
14
15 k-mer spectra: two well-used metrics in metagenomics (Bray-Curtis and Jensen distance), as
16
17 well as a cosine similarity metric. Cosine similarity, although extensively used in computer
18
19 science, has been rarely implemented in genomic and metagenomic studies [46]. To our
20
21 knowledge, this work is the first to describe the use of the cosine similarity metric to cluster
22
23 metagenomes based on their k-mer content.
24
25
26
27
28

29 Finally, the analysis of large-scale metagenomic analysis requires access to large computing
30
31 resources. In order to use Libra, the user requires access to a Hadoop framework. In order
32
33 allow for a better access to the tool and to computing resources, we provide a web-based
34
35 implementation tool embedded in the CyVerse advanced cyberinfrastructure through iMicrobe
36
37 [37]. The work described here is the first step in implementing a free cloud-based computing
38
39 resource for *de-novo* comparative metagenomics that can be broadly used by scientists to
40
41 analyze large-scale shared data resources. Moreover, the code can be ported to any Hadoop
42
43 cluster (e.g., Wrangler at TACC, Amazon EMR, or private Hadoop clusters). This computing
44
45 paradigm is consistent with recent efforts to increase the accessibility of big data sets in the
46
47 cloud, such as the Pan-Cancer Analyses of Whole Genomes Project [56].
48
49
50
51
52

53 **METHODS**

54 **Libra Algorithm Detailed Description:**

55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 **k-mer size.** Libra calculates the distances between samples based on their k-mer composition.

5
6 The canonical representation of the k-mer is used to reduce the number of stored k-mers.

7
8 Several considerations should be taken into account for choosing the k-mer size k . Larger

9
10 values of k result in fewer matches due to sequencing errors and fragmentary metagenomic

11
12 data. However, smaller values of k give less information about the sequence similarities. In

13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
Libra, k is a configurable parameter chosen by the user and is set by default to k equal to 21.

This value was reported to be at the inflection point where the k-mer matches move from random to a representative of the read content and is generally resilient to sequencing error and variation [57,58].

Distance Matrix Computation. Libra provides three distance metrics — Cosine Similarity, Bray-Curtis, and Jensen-Shannon. Cosine Similarity is the default.

Cosine Similarity Metric. Libra constructs a vector v_s for each sample s from the weight of each k-mer k in the sample ($w_{k,s}$). Each dimension in the vector corresponds to the weight of the corresponding k-mer:

$$v_s = (w_{k1,s}, w_{k2,s}, w_{k3,s}, \dots, w_{kn,s})$$

The weight of a k-mer in a sample (w_{ks}) can be derived from the frequency of the k-mer (f_{ks}) in several ways. The simplest uses the raw frequency of the k-mer ($w_{ks} = f_{ks}$), called *Natural Weighting*. Another uses *Logarithmic Weighting* ($w_{ks} = 1 + \log(f_{ks})$) to not give too much weight to highly abundant k-mers. In this weighting w_{ks} grows logarithmically with the frequency f_{ks} , reducing the effect on the distance of highly abundant k-mers caused by sequencing artifacts.

Once their vectors have been constructed, the distance between two samples (s_1 and s_2) is derived using distance metrics. For example, the distance between the two samples using Cosine Similarity is determined as follows:

$$\begin{aligned} \text{Distance}(s_1, s_2) &= 1 - \text{CosineSimilarity}(s_1, s_2) \\ &= 1 - \cos(v_{s_1}, v_{s_2}) = 1 - \frac{v_{s_1} \cdot v_{s_2}}{\|v_{s_1}\| \times \|v_{s_2}\|} = 1 - \frac{D_{s_1, s_2}}{M_{s_1} \times M_{s_2}} \end{aligned}$$

$$\text{where, } D_{s_1, s_2} = v_{s_1} \cdot v_{s_2} = \sum_{i \in s_1 \cap s_2} w_{ki, s_1} \times w_{ki, s_2},$$

$$M_s = \|v_s\| = \sqrt{\sum_{i \in s} (w_{ki, s})^2}$$

In other words, D_{s_1, s_2} is the dot product of the vectors v_{s_1} and v_{s_2} , and M_s is the magnitude (length) of the vector v_s . The distance between two NGS samples is the cosine of the angle between their vectors v_s ; the magnitude of the vector M_s is not taken into account in the metric thereby normalizing samples with different numbers of total base pairs.

Inverted Index Construction. A naïve implementation would require the storage of one vector with 4^k dimensions per sample, where k is the k-mer length. For a k of 21, each vector would have more than one million dimensions. To reduce the overhead, Libra stores and computes the distance on a single *inverted index* with the k-mer frequencies from multiple samples and performs the distance computation on the index directly. The inverted index is indexed by k-mer, and each entry is an index record containing a list of pairs, each of which contains a sample identifier and the frequency of the k-mer in the sample.

$$\text{index record} = k\text{-mer} : \{ \langle \text{sample-id}, \text{frequency} \rangle, \langle \text{sample-id}, \text{frequency} \rangle \dots \}$$

The records in the index are stored in an alphabetical order by k-mer, allowing the record for a particular k-mer to be found via binary search. The k-mer record contains the k-mer frequency in

1
2
3
4 each sample, not the weight, to allow for different weighting functions to be applied during
5
6 distance matrix computation.
7
8

9
10 **Sweep line algorithm.** To compute the distance between two samples S_1 and S_2 , Libra must
11
12 compute the three values D_{s_1,s_2} , M_{s_1} , and M_{s_2} . The values are calculated by scanning through
13
14 the vectors v_{s_1} and v_{s_2} and computing the values. The time for the distance matrix computation
15
16 is proportional to the number of dimensions (the number of k-mers) in the two vectors. In
17
18 general, computing all-vs-all comparisons on n samples would require $n \times (n - 1) / 2$ vector scans,
19
20 which becomes prohibitively expensive as n gets large. Libra uses a sweep line algorithm [38] to
21
22 greatly reduce the computational time. The sweep line algorithm only requires a single scan of
23
24 all vectors to compute the distance of all pairs of samples ([See Supplemental Figure 4]). Briefly,
25
26 Libra sweeps a line through all the vectors simultaneously starting with the first component.
27
28

29
30 Libra outputs a record of the non-zero values of the following format:
31
32

$$33 \quad \text{record} = k\text{-mer} : \{ \langle \text{sample-id}, \text{weight} \rangle, \langle \text{sample-id}, \text{weight} \rangle, \dots \}$$

34
35
36
37 Libra then moves the sweep line to the next component and performs the same operation. From
38
39 the output records, contributions to M_s for each sample in the record are computed and
40
41 accumulated. Contributions to D are also computed from the record by extracting sample pairs.
42
43

44 For example, the record $\{ \langle s_1, x \rangle, \langle s_2, y \rangle, \langle s_4, z \rangle \}$ has three sample pairs

45
46 (s_1s_2) , (s_1s_4) and (s_2s_4) . Libra then computes contribution to D for each pair, e.g. $x * y$ is added
47
48 to D_{s_1,s_2} , $x * z$ is added to D_{s_1,s_4} , and $y * z$ is added to D_{s_2,s_4} . Using this method, Libra
49
50

51
52 computes the distances of every sample pairs in an input dataset in linear time. Other distance
53
54 metrics, such as Bray-Curtis and Jensen-Shannon, can also be computed in the same fashion.
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 The sweep algorithm is particularly easy to implement on an inverted index; it consists of simply
5
6 stepping through the (sorted) k-mers. Furthermore, the sweep algorithm is easily parallelized.
7

8
9 The k-mer space is partitioned and a separate sweep is performed on each partition computing
10
11 the contributions of its k-mer frequencies to the D and M values. At the end of the
12
13 computation, the intermediate D and M values are combined together to produce the final D
14
15 and M values and thereby the distance matrix. Each sweep uses binary search to find the first
16
17 k-mer in the partition.
18
19
20

21 **Terabyte Sort.** Libra groups the samples automatically based on the number and size (by
22
23 default 4GB per group). Similar to Terabyte Sort [59] the index records are partitioned by k-mer
24
25 ranges and the records in each partition is stored in a separate *chunk file*. All k-mers in partition
26
27 n appear before the k-mers in partition $n + 1$ in lexicographic order. This facilitates breaking
28
29 computation and I/O down into smaller tasks, so that work of creating an index can be
30
31 distributed across several machines.
32
33
34
35

36 **k-mer space partitioning.** Both the inverted index construction and the distance matrix
37
38 computation require partitioning the k-mer space so that different partitions can be processed
39
40 independently. For the partitioning to be effective, the workload should be balanced across the
41
42 partitions. Simply partitioning into fixed-size partitions based on the k-mer space will not ensure
43
44 balanced workloads, as the k-mers do not appear with uniform frequency. Some partitions may
45
46 have more k-mer records than others, and thereby incur higher processing costs. Instead, the
47
48 partitions should be created based on the k-mer distribution, so that each partition has roughly
49
50 the same number of records ([See Supplemental Figure 5]).
51
52
53
54

55
56 Computing the exact k-mer distribution across all the samples is too expensive in both space
57
58 and time, therefore Libra approximates the distribution instead. A histogram is constructed using
59
60
61
62
63
64
65

1
2
3
4 the first 6 letters of the k-mers in each sample, which requires much less space and time to
5
6 compute. In practice, partitioning based on this histogram adequately partitions the k-mer space
7
8 so that the workloads are sufficiently balanced across the partitions.
9

10
11 **Scalability benchmarking for Libra.** We used synthetic datasets for a scalability benchmark.
12
13 Each dataset contains 10 billion bytes (approximately 9.3 GB). We used four datasets
14
15 consisting of 10 (93GB), 20 (186GB), 30 (279GB) and 40 (372GB) samples in the benchmark.
16
17 Each experiment was run three times, and an average of the three runs reported ([See
18
19 Supplemental Table 4 for details]). The runtime of Libra increased linearly with increased input
20
21 volume (Figure 6). This shows that Libra efficiently handles the increased volume of input and
22
23 efficiently computes distances between all sample pairs while the number of sample pairs
24
25 increases quadratically.
26
27
28
29
30

31 **Figure 6 Scalability testing for Libra.** Runtimes of Libra on four datasets consisting of 10, 20,
32
33 30 and 40 samples (total sizes of 93GB, 186GB, 279GB, and 372GB, respectively). Libra was
34
35 performed with default parameters (k=20, Logarithmic weighting, and Cosine Similarity).
36
37 Runtimes were averaged out over 3 runs. The total runtime of Libra increased linearly with
38
39 increased input volume. Both index construction and distance matrix computation showed
40
41 linearly increased runtimes for the increased input volume. This shows that Libra performs
42
43 efficiently and scales to input although the number of distances between sample pairs to be
44
45 computed increases quadratically.
46
47
48
49
50

51
52 **Benchmarking runtimes of different distance metrics in Libra.** We used the same synthetic
53
54 dataset with 40 samples (372GB in total) in the scalability benchmarking (Figure 7). We
55
56 measured the runtimes of Libra for the different distance metrics. Once the index is constructed
57
58
59
60
61
62
63
64
65

1
2
3
4 all distance metrics are calculated using that index; thus, runtimes of the inverted index
5
6 construction for the different metrics are the same. Each experiment was run three times and
7
8 the average reported ([See Supplemental Table 4 for details]). Differences in runtimes are
9
10 mainly due to the different computational workload of distance metrics (Figure 7). For example,
11
12 Jensen-Shannon requires more multiplications and divisions in nested loops than Cosine
13
14 Similarity, incurring more computational workload. Yet, distance matrix computation with
15
16 Jensen-Shannon took only 12.64% of total runtime.
17
18
19
20

21 **Figure 7. Runtime for different distance metrics.** Runtimes for three different distance
22
23 metrics (Cosine Similarity, Bray-Curtis, and Jensen-Shannon) in Libra with 40 samples of input
24
25 (372GB in total). Libra was performed with default parameters (k=20 and Logarithmic
26
27 weighting). Runtimes were averaged over 3 runs. An inverted index was reused for all three
28
29 distance metrics because the inverted index Libra constructs are independent of the distance
30
31 metrics. Cosine Similarity took the shortest runtime among the three metrics while
32
33 Jensen-Shannon took the longest. Jensen-Shannon took almost twice as long as Cosine
34
35 Similarity because it requires more mathematical computations. Because of its fastest runtime,
36
37 Cosine Similarity is used by default in Libra.
38
39
40
41

42 **Advanced cyberinfrastructure for Libra in iMicrobe.** To improve access to Libra we made it
43
44 available on the iMicrobe website [37]. A researcher with a CyVerse account can run Libra on
45
46 iMicrobe by filling out a simple web form specifying the input files and parameters. Input files are
47
48 selected from the CyVerse Data Store where they have either been uploaded by the user to
49
50 their home directory or are part of the iMicrobe Data Commons. When a job is submitted, the
51
52 user is presented with the status of the job, and on completion the output files and visualization
53
54 of results. To deploy Libra on iMicrobe, we developed a job dispatch service to automate the
55
56 execution of Libra on a University of Arizona Hadoop cluster. The service is written in NodeJS
57
58
59
60
61
62
63
64
65

1
2
3
4 and accepts a JSON description of the job inputs and parameters, stages the input files onto the
5
6 UA Hadoop cluster, executes Libra with the given parameters, and transfers the resulting output
7
8 files to the user's home directory in the CyVerse Data Store. The service provides a RESTful
9
10 interface that mimics the Agave API Jobs service and is secured using an Agave OAuth2 token.
11
12
13 The source code is available on Github [60].
14

15 16 **Experimental Environment Description:**

17
18 **Mash and Simka configurations.** Mash v1.1 was run on the metagenomic datasets with the
19
20 following parameters: `-r -s 10000 -m 2` [19]. The analysis of assemblies was run without the
21
22 parameter “-r”, used for short sequences.
23

24
25 Simka v1.3.2 was run on the metagenomic datasets with the following parameters:

26
27 `-abundance-min 2 -max-reads [MINCOUNT] -simple-dist -complex-dist`, where [MINCOUNT] is
28
29 the smallest sequence count across the analyzed samples.
30

31
32
33 **Hadoop cluster configuration.** The Libra experiments described in the paper were performed
34
35 on a Hadoop cluster consisting of 10 physical nodes (9 MapReduce worker nodes). Each node
36
37 contains 12 CPUs and 128 GB of RAM and is configured to run a maximum of 7 YARN
38
39 containers simultaneously with 10 GB of RAM per container. The remaining system resources
40
41 are reserved for the operating system and other Hadoop services such as Hive or HBase.
42
43
44

45
46 **The rationale for not porting Libra to Spark.** Spark [61] is increasingly popular for scientific
47
48 data analysis [62] because of its outstanding performance provided by fast in-memory
49
50 processing. Although Libra is currently implemented on Hadoop MapReduce, Libra can be
51
52 easily ported to Spark because both Hadoop MapReduce and Spark have similar interfaces for
53
54 data processing and partitioning. For example, Resilient Distributed Datasets (RDD) can be
55
56 partitioned and distributed over a Spark cluster using Libra's k-mer range partitioning. RDDs are
57
58
59
60
61
62
63
64
65

1
2
3
4 memory-resident, allowing Spark to significantly improve the performance of Libra's k-mer
5
6 counting and distance matrix computation by avoiding slow disk I/O for intermediate data. We
7
8 implemented Libra using Hadoop MapReduce because Spark requires much more RAM than
9
10 Hadoop MapReduce, significantly increasing the cost of the cluster.
11
12

13 **AVAILABILITY AND IMPLEMENTATION**

14
15
16 **Project home page:** Program binary, source code and documentation for Libra are available in
17
18 Github [63]; Libra web-based App is in iMicrobe [37] under Apps; code to implement the Libra
19
20 web-based App is in Github [60]; **Operating system(s):** MapReduce 2.0 (Apache Hadoop 2.3.0
21
22 or above); **Programming language:** Java 7 (or above); **Other requirements:** none; **License:**
23
24 Apache License Version 2.0; **Any restrictions to use by non-academics:** no license needed.
25
26

27
28 Libra has been registered with the SciCrunch database under reference ID: SCR_016608.
29

30 **AVAILABILITY OF SUPPORTING DATA**

31
32
33 Snapshots of the code and other supporting data are available in the GigaScience repository,
34
35 GigaDB [64].
36

37 **ABBREVIATIONS**

38
39 HDFS - high-performance distributed file system; HPC- high-performance computer cluster; GB
40
41 - gigabytes; TOV - Tara Ocean Viromes; HMP - Human Microbiome Project; GOS - Global
42
43 Ocean Survey; ENA - European Nucleotide Archive; CAMI - Critical Assessment of
44
45 Metagenome Interpretation
46
47
48
49

50 **COMPETING INTERESTS**

51
52
53 The authors declare no competing interests.
54

55 **FUNDING**

56
57
58
59
60
61
62
63
64
65

1
2
3
4 This work was supported by the National Science Foundation award #1640775 to BLH and
5
6 JHH. System support and access for the Hadoop cluster was provided by University of Arizona
7
8 Information Technology Services. The following reagent was obtained through BEI Resources,
9
10 NIAID, NIH as part of the Human Microbiome Project: Genomic DNA from Microbial Mock
11
12 Community B (Staggered, High Concentration), v5.2H, for Whole Genome Shotgun
13
14 Sequencing, HM-277D.
15

16 17 18 **AUTHORS' CONTRIBUTIONS**

19
20 BLH, JHH, and IC conceived of the Libra algorithm and code. IC wrote the Libra code. BLH,
21
22 AJP, and IC planned and carried out the analyses. BLH, AJP, JHH, and IC contributed to the
23
24 interpretation of the results. MB and KYC implemented Libra as an application in iMicrobe. BLH,
25
26 AJP, JHH, and IC wrote the manuscript. All authors provided critical feedback and helped shape
27
28 the research, analysis, and manuscript.
29
30

31 32 **ACKNOWLEDGMENTS**

33
34 We thank Binil Benjamin, Russell Lewis and members of the Hurwitz Lab for comments on the
35
36 Libra algorithm and feedback on the manuscript. We thank the staff at the University of Arizona
37
38 Information Technology Services for access and system support for the Hadoop cluster. We
39
40 thank George Watts for feedback and critical discussions on the manuscript and support for
41
42 sequencing the mock community.
43
44

45 46 **REFERENCES**

- 47
48
49 1. Yooseph S, Sutton G, Rusch DB, Halpern AL, Williamson SJ, Remington K, et al. The
50 Sorcerer II Global Ocean Sampling expedition: expanding the universe of protein families. PLoS
51 Biol. 2007;5:e16.
52
53 2. Sunagawa S, Coelho LP, Chaffron S, Kultima JR, Labadie K, Salazar G, et al. Structure and
54 function of the global ocean microbiome. Science [Internet]. sciencemag.org; 2015;348.
55 Available from: <http://www.sciencemag.org/content/348/6237/1261359.abstract>
56
57 3. Hurwitz BL, Sullivan MB. The Pacific Ocean Virome (POV): a marine viral metagenomic
58 dataset and associated protein clusters for quantitative viral ecology. PLoS One.
59
60
61
62
63
64
65

1
2
3
4 2013;8:e57355.

5
6 4. Marcy Y, Ouverney C, Bik EM, Lösekann T, Ivanova N, Martin HG, et al. Dissecting biological
7 “dark matter” with single-cell genetic analysis of rare and uncultivated TM7 microbes from the
8 human mouth. *Proc Natl Acad Sci U S A. National Academy of Sciences*; 2007;104:11889–94.
9

10
11 5. Rinke C, Schwientek P, Sczyrba A, Ivanova NN, Anderson IJ, Cheng J-F, et al. Insights into
12 the phylogeny and coding potential of microbial dark matter. *Nature. nature.com*;
13 2013;499:431–7.
14

15
16 6. Hug LA, Baker BJ, Anantharaman K, Brown CT, Probst AJ, Castelle CJ, et al. A new view of
17 the tree of life. *Nat Microbiol. nature.com*; 2016;1:16048.
18

19
20 7. Dubinkina VB, Ischenko DS, Ulyantsev VI, Tyakht AV, Alexeev DG. Assessment of k-mer
21 spectrum applicability for metagenomic dissimilarity analysis. *BMC Bioinformatics*.
22 *bmcbioinformatics.biomedcentral. ...*; 2016;17:38.
23

24
25 8. Teeling H, Waldmann J, Lombardot T, Bauer M, Glockner FO. TETRA: a web-service and a
26 stand-alone program for the analysis and comparison of tetranucleotide usage patterns in DNA
27 sequences. *BMC Bioinformatics*. 2004;5:163.
28

29
30 9. Wu Y-W, Ye Y. A novel abundance-based algorithm for binning metagenomic sequences
31 using l-tuples. *J Comput Biol. online.liebertpub.com*; 2011;18:523–34.
32

33
34 10. Fofanov Y, Luo Y, Katili C, Wang J, Belosludtsev Y, Powdrill T, et al. How independent are
35 the appearances of n-mers in different genomes? *Bioinformatics. Oxford University Press*;
36 2004;20:2421–8.
37

38
39 11. Maillet N, Lemaitre C, Chikhi R, Lavenier D, Peterlongo P. Compareads: comparing huge
40 metagenomic experiments. *BMC Bioinformatics. bmcbioinformatics.biomedcentral. ...*; 2012;13
41 Suppl 19:S10.
42

43
44 12. Maillet N, Collet G, Vannier T, Lavenier D, Peterlongo P. Commet: Comparing and
45 combining multiple metagenomic datasets. 2014 IEEE International Conference on
46 Bioinformatics and Biomedicine (BIBM). *ieeexplore.ieee.org*; 2014. p. 94–8.
47

48
49 13. Ulyantsev VI, Kazakov SV, Dubinkina VB, Tyakht AV, Alexeev DG. MetaFast: fast
50 reference-free graph-based comparison of shotgun metagenomic data. *Bioinformatics. Oxford*
51 *Univ Press*; 2016;32:2760–7.
52

53
54 14. Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, et al. Mash: fast
55 genome and metagenome distance estimation using MinHash. *Genome Biol. biorxiv.org*;
56 2016;17:132.
57

58
59 15. Benoit G, Peterlongo P, Mariadassou M, Drezen E, Schbath S, Lavenier D, et al. Multiple
60 comparative metagenomics using multiset k-mer counting. *PeerJ Comput Sci. PeerJ Inc.*;
61 2016;2:e94.
62

63
64 16. Broder AZ. On the resemblance and containment of documents. *Proceedings Compression*
65 *and Complexity of SEQUENCES 1997 (Cat No97TB100171) [Internet]*. Available from:

1
2
3
4 <http://dx.doi.org/10.1109/sequen.1997.666900>

5
6
7 17. Koslicki D, Zabeti H. Improving Min Hash via the Containment Index with applications to
8 Metagenomic Analysis [Internet]. bioRxiv. 2017 [cited 2018 Oct 19]. p. 184150. Available from:
9 <https://www.biorxiv.org/content/early/2017/09/04/184150.abstract>

10
11 18. Brown CT, Irber L. sourmash: a library for MinHash sketching of DNA. The Journal of Open
12 Source Software [Internet]. theoj.org; 2016;1. Available from:
13 <http://www.theoj.org/joss-papers/joss.00027/10.21105.joss.00027.pdf>

14
15 19. Seth S, Välimäki N, Kaski S, Honkela A. Exploration and retrieval of whole-metagenome
16 sequencing samples. Bioinformatics. academic.oup.com; 2014;30:2471–9.

17
18
19 20. Dean J, Ghemawat S. MapReduce: Simplified Data Processing on Large Clusters. Commun
20 ACM. New York, NY, USA: ACM; 2008;51:107–13.

21
22 21. Kolker N, Higdon R, Broomall W, Stanberry L, Welch D, Lu W, et al. Classifying proteins into
23 functional groups based on all-versus-all BLAST of 10 million proteins. OMICS.
24 online.liebertpub.com; 2011;15:513–21.

25
26 22. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The
27 Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA
28 sequencing data. Genome Res. genome.cshlp.org; 2010;20:1297–303.

29
30
31 23. Langmead B, Schatz MC, Lin J, Pop M, Salzberg SL. Searching for SNPs with cloud
32 computing. Genome Biol. biomedcentral.com; 2009;10:R134.

33
34 24. Wall DP, Kudtarkar P, Fusaro VA, Pivovarov R, Patil P, Tonellato PJ. Cloud computing for
35 comparative genomics. BMC Bioinformatics. 2010;11:259.

36
37
38 25. Langmead B, Hansen KD, Leek JT. Cloud-scale RNA-sequencing differential expression
39 analysis with Myrna. Genome Biol. genomebiology.biomedcentral.com; 2010;11:R83.

40
41 26. Jourden L, Bernard M, Dillies M-A, Le Crom S. Eoulsan: a cloud computing-based
42 framework facilitating high throughput sequencing analyses. Bioinformatics. academic.oup.com;
43 2012;28:1542–3.

44
45 27. Nguyen T, Shi W, Ruden D. CloudAligner: A fast and full-featured MapReduce based tool
46 for sequence mapping. BMC Res Notes. biomedcentral.com; 2011;4:171.

47
48
49 28. Schatz MC. BlastReduce: high performance short read mapping with MapReduce.
50 University of Maryland, <http://cgis.cs.umd.edu/Grad/scholarlypapers/papers/MichaelSchatz.pdf>
51 [Internet]. cs.umd.edu; 2008; Available from:
52 https://www.cs.umd.edu/sites/default/files/scholarly_papers/MichaelSchatz_1.pdf

53
54
55 29. Schatz MC. CloudBurst: highly sensitive read mapping with MapReduce. Bioinformatics.
56 Oxford Univ Press; 2009;25:1363–9.

57
58 30. Pandey RV, Schlötterer C. DistMap: a toolkit for distributed short read mapping on a
59 Hadoop cluster. PLoS One. dx.plos.org; 2013;8:e72614.

60
61
62
63
64
65

31. Nordberg H, Bhatia K, Wang K, Wang Z. BioPig: a Hadoop-based analytic toolkit for large-scale sequence data. *Bioinformatics*. Oxford Univ Press; 2013;29:3014–9.
32. Gao T, Guo Y, Wei Y, Wang B, Lu Y, Cicotti P, et al. Bloomfish: A Highly Scalable Distributed K-mer Counting Framework. 2017 IEEE 23rd International Conference on Parallel and Distributed Systems (ICPADS). ieeexplore.ieee.org; 2017. p. 170–9.
33. Menon RK, Bhat GP, Schatz MC. Rapid Parallel Genome Indexing with MapReduce. *Proceedings of the Second International Workshop on MapReduce and Its Applications*. New York, NY, USA: ACM; 2011. p. 51–8.
34. Michie MG. Use of the Bray-Curtis similarity measure in cluster analysis of foraminiferal data. *Math Geol*. Kluwer Academic Publishers-Plenum Publishers; 1982;14:661–7.
35. Lin J. Divergence measures based on the Shannon entropy. *IEEE Trans Inf Theory*. 1991;37:145–51.
36. Brum JR, Ignacio-Espinoza JC, Roux S, Doucier G, Acinas SG, Alberti A, et al. Patterns and ecological drivers of ocean viral communities. *Science* [Internet]. [sciencemag.org](http://www.sciencemag.org); 2015;348. Available from: <http://www.sciencemag.org/content/348/6237/1261498.abstract>
37. Ken Youens-Clark, Matthew Bomhoff, and Bonnie L Hurwitz. iMicrobe [Internet]. iMicrobe. 2018 [cited 2018 Dec 7]. Available from: <http://imicrobe.us>
38. Goff SA, Vaughn M, McKay S, Lyons E, Stapleton AE, Gessler D, et al. The iPlant Collaborative: Cyberinfrastructure for Plant Biology. *Front Plant Sci*. ncbi.nlm.nih.gov; 2011;2:34.
39. Devisetty UK, Kennedy K, Sarando P, Merchant N, Lyons E. Bringing your tools to CyVerse Discovery Environment using Docker. *F1000Res*. 2016;5:1442.
40. Alise Ponsero and Bonnie L. Hurwitz. Simulated Metagenomics Mock Communities [Internet]. *Libra: scalable k-mer based tool for massive all-vs-all metagenome comparisons*. 2018 [cited 2018 Dec 11]. Available from: <https://doi.org/10.7946/MQ0G>
41. McElroy KE, Luciani F, Thomas T. GemSIM: general, error-model based simulator of next-generation sequencing data. *BMC Genomics*. biomedcentral.com; 2012;13:74.
42. Stöcker BK, Köster J, Rahmann S. SimLoRD: Simulation of Long Read Data. *Bioinformatics*. academic.oup.com; 2016;32:2704–6.
43. Human Microbiome Project Consortium. Structure, function and diversity of the healthy human microbiome. *Nature*. nature.com; 2012;486:207–14.
44. Diepenbroek M, Grobe H, Reinke M, Schindler U, Schlitzer R, Sieger R, et al. PANGAEA—an information system for environmental sciences. *Comput Geosci*. Elsevier; 2002;28:1201–10.
45. Alice McHardy, and Alexander Sczyrba. Critical Assessment of Metagenome Interpretation (CAMI) 2nd Challenge [Internet]. CAMI. 2018 [cited 2018 Dec 11]. Available from:

1
2
3
4 <https://data.cami-challenge.org/participate>

5
6 46. Okuda S, Tsuchiya Y, Kiriya C, Itoh M, Morisaki H. Virtual metagenome reconstruction
7 from 16S rRNA gene sequences. *Nat Commun.* nature.com; 2012;3:1203.

8
9 47. Watts GS, Youens-Clark K, Slepian MJ, Wolk DM, Oshiro MM, Metzger GS, et al. 16S rRNA
10 gene sequencing on a benchtop sequencer: accuracy for identification of clinically important
11 bacteria. *J Appl Microbiol.* Wiley Online Library; 2017;123:1584–96.

12
13 48. Bergh O, Borsheim KY, Bratbak G, Heldal M. High abundance of viruses found in aquatic
14 environments. *Nature.* 1989;340:467–8.

15
16 49. Hurwitz BL, Brum JR, Sullivan MB. Depth-stratified functional and taxonomic niche
17 specialization in the “core” and “flexible” Pacific Ocean Virome. *ISME J [Internet].* nature.com;
18 2014; Available from:
19 <http://www.nature.com/ismej/journal/vaop/ncurrent/full/ismej2014143a.html>

20
21 50. Minot S, Wu GD, Lewis JD, Bushman FD. Conservation of gene cassettes among diverse
22 viruses of the human gut. *PLoS One.* 2012;7:e42342.

23
24 51. Sun Y, Cai Y, Huse SM, Knight R, Farmerie WG, Wang X, et al. A large-scale benchmark
25 study of existing algorithms for taxonomy-independent microbial community analysis. *Brief*
26 *Bioinform.* Oxford Univ Press; 2012;13:107–21.

27
28 52. Edgar RC. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics.*
29 2010;26:2460–1.

30
31 53. Niu BF, Fu LM, Sun SL, Li WZ. Artificial and natural duplicates in pyrosequencing reads of
32 metagenomic data. *BMC Bioinformatics.* 2010;11:187.

33
34 54. Cai Y, Sun Y. ESPRIT-Tree: hierarchical clustering analysis of millions of 16S rRNA
35 pyrosequences in quasilinear computational time. *Nucleic Acids Res.* Oxford Univ Press;
36 2011;39:e95.

37
38 55. Hurwitz BL, Brum JR, Sullivan MB. Depth-stratified functional and taxonomic niche
39 specialization in the “core” and “flexible” Pacific Ocean Virome. *ISME J.* 2015;9:472–84.

40
41 56. Weinstein JN, Collisson EA, Mills GB, Shaw KRM, Ozenberger BA, Ellrott K, et al. The
42 cancer genome atlas pan-cancer analysis project. *Nat Genet.* Nature Publishing Group;
43 2013;45:1113–20.

44
45 57. Hurwitz BL, Westveld AH, Brum JR, Sullivan MB. Modeling ecological drivers in marine viral
46 communities using comparative metagenomics and network analyses. *PNAS.*
47 2014;111:10714–9.

48
49 58. Kurtz S, Narechania A, Stein JC, Ware D. A new method to compute K-mer frequencies and
50 its application to annotate large repetitive plant genomes. *BMC Genomics.* 2008;9:517.

51
52 59. O'Malley O. Terabyte sort on apache hadoop. Yahoo, available online at:
53 <http://sortbenchmark.org/Yahoo-Hadoop.pdf>, (May). Citeseer; 2008;1–3.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

60. Bomhoff M. Ocean Cloud Commons Plan B Github [Internet]. OCC Plan B Github. 2018 [cited 2018 Dec 11]. Available from: <https://github.com/hurwitzlab/occ-plan-b>

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

62. Guo R, Zhao Y, Zou Q, Fang X, Peng S. Bioinformatics applications on Apache Spark. Gigascience [Internet]. academic.oup.com; 2018; Available from: <http://dx.doi.org/10.1093/gigascience/giy098>

63. Choi I. Libra Github [Internet]. Libra Github. 2018 [cited 2018 Dec 11]. Available from: <https://www.github.com/iychoi/Libra>

64. Choi I; Ponsero AJ; Bomhoff M; Youens-Clark K; Hartman JH; Hurwitz BL: Supporting data for "Libra: scalable k-mer based tool for massive all-vs-all metagenome comparisons." GigaScience Database. 2018. <http://dx.doi.org/10.5524/100547>.

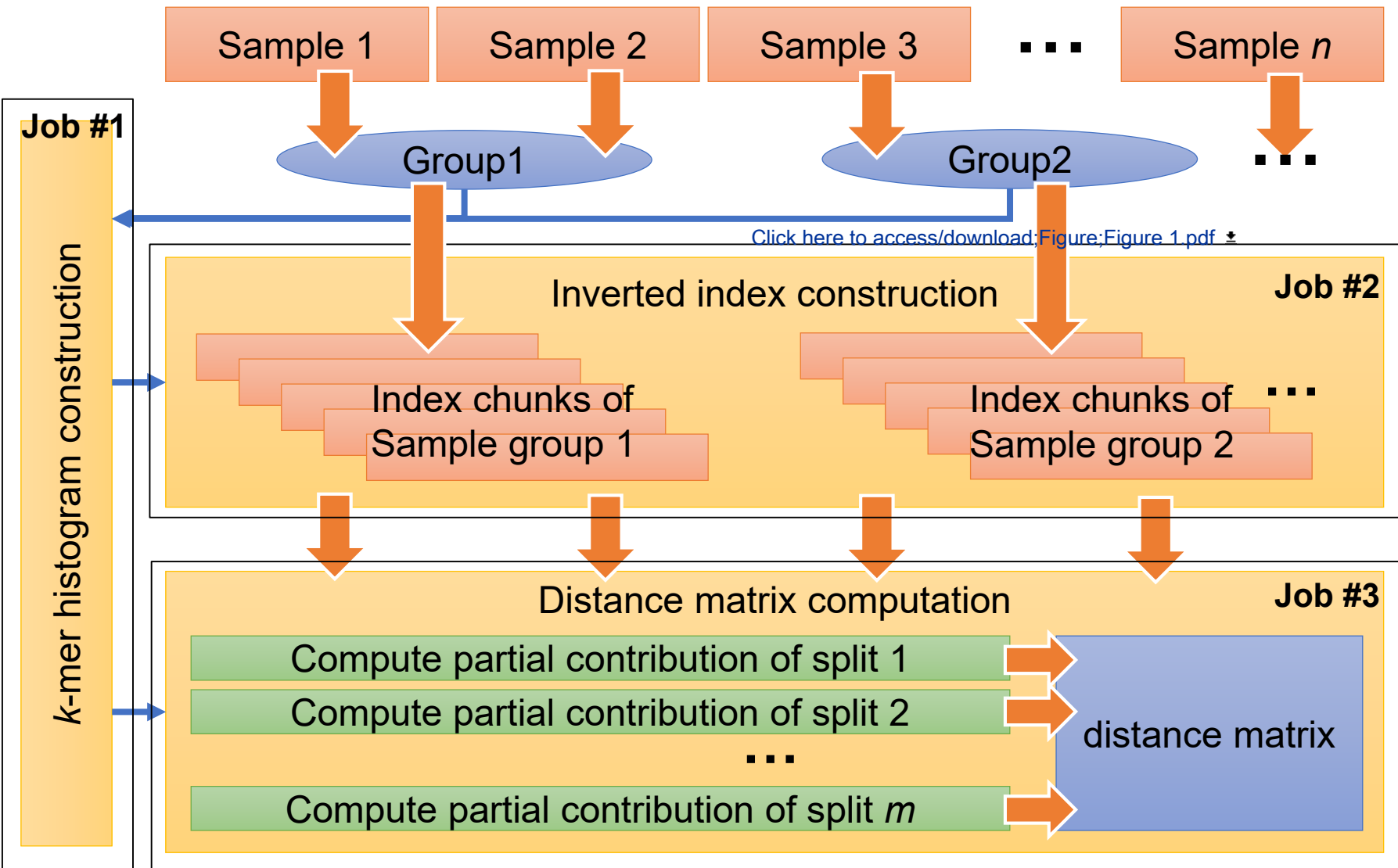
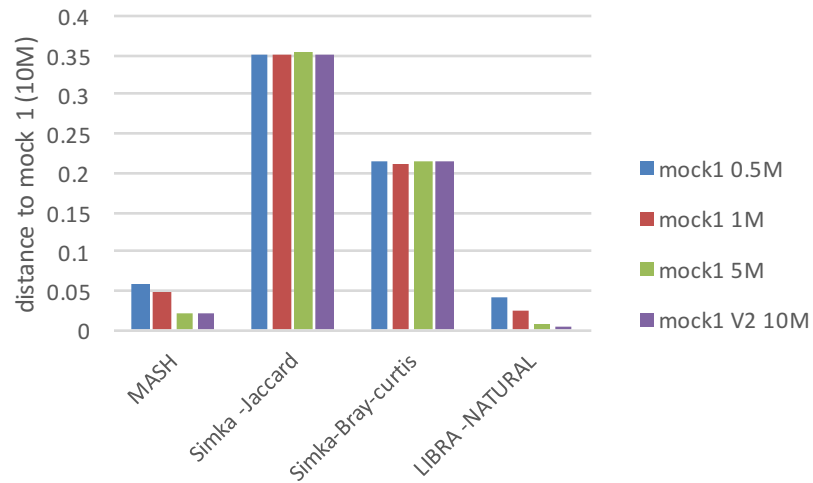


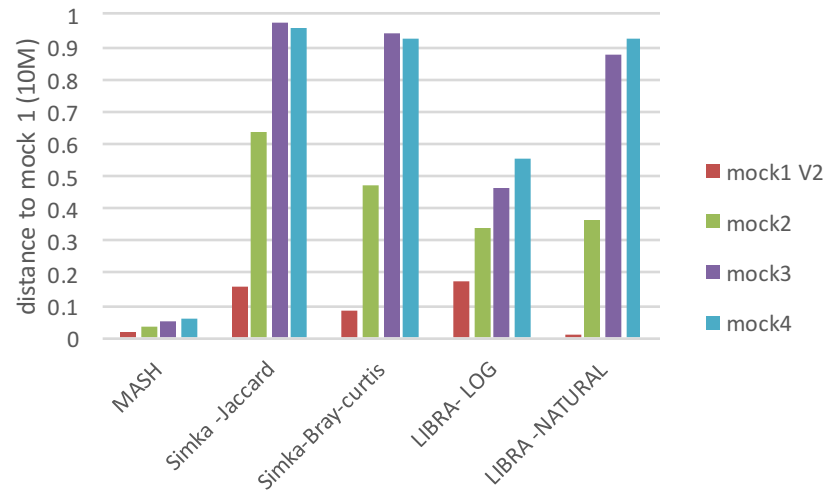
Figure 2

[Click here to access/download;Figure;Figure 2.pdf](#)

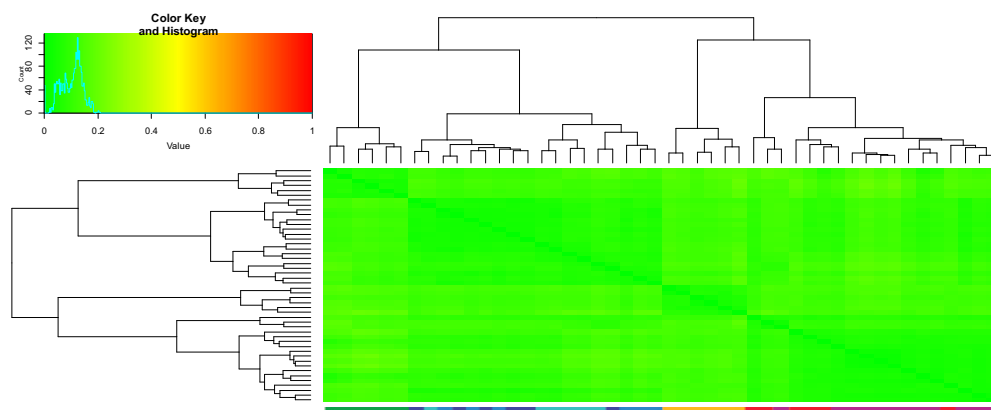
A



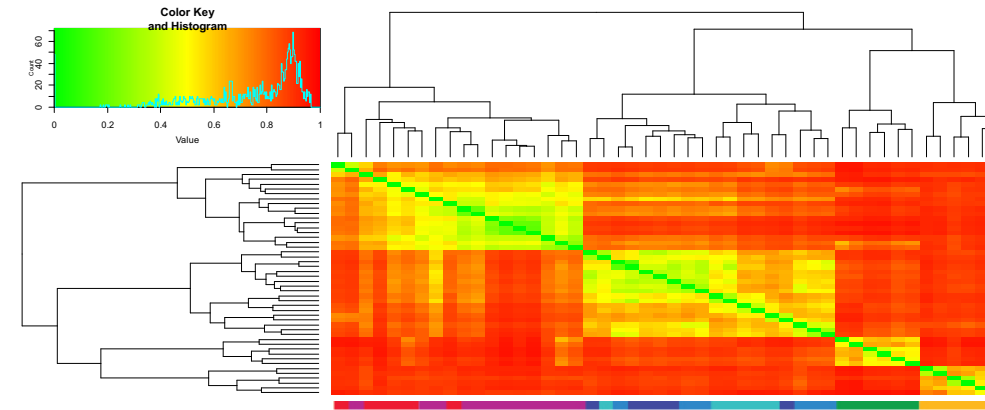
B



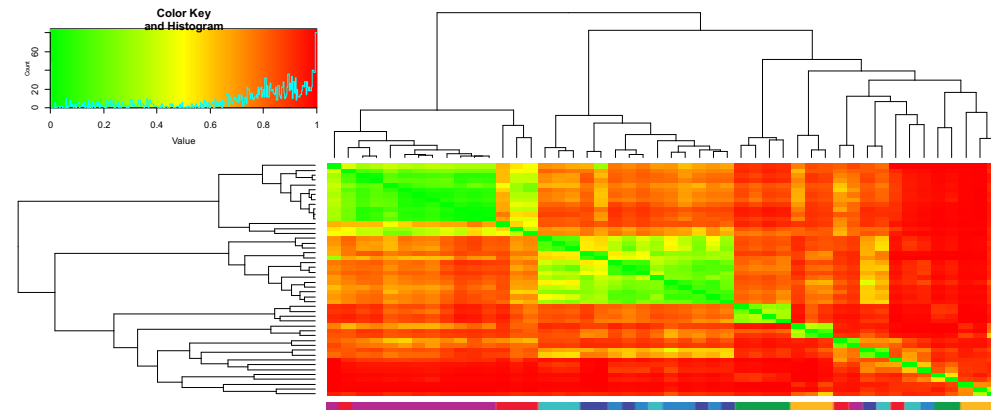
A - MASH



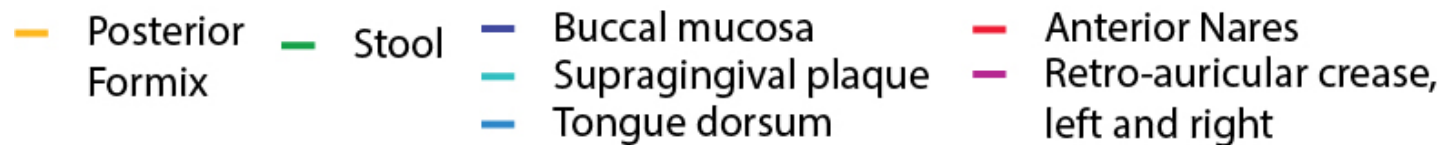
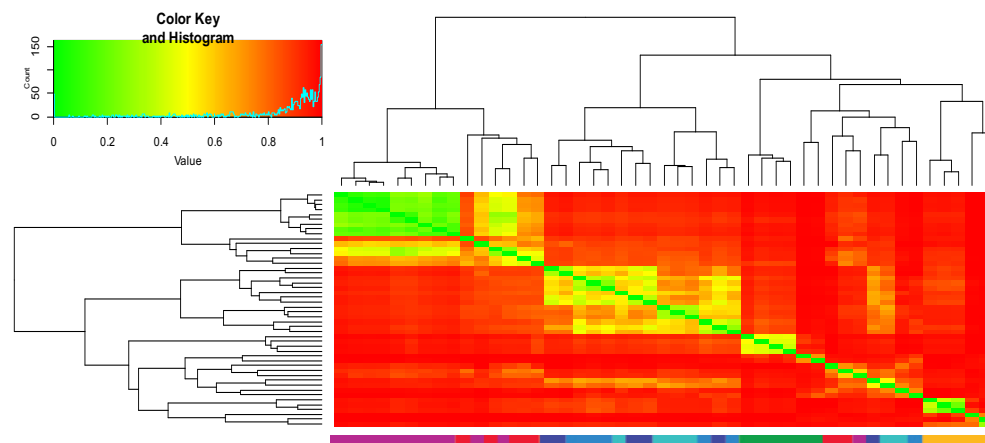
B – LIBRA, log weighting



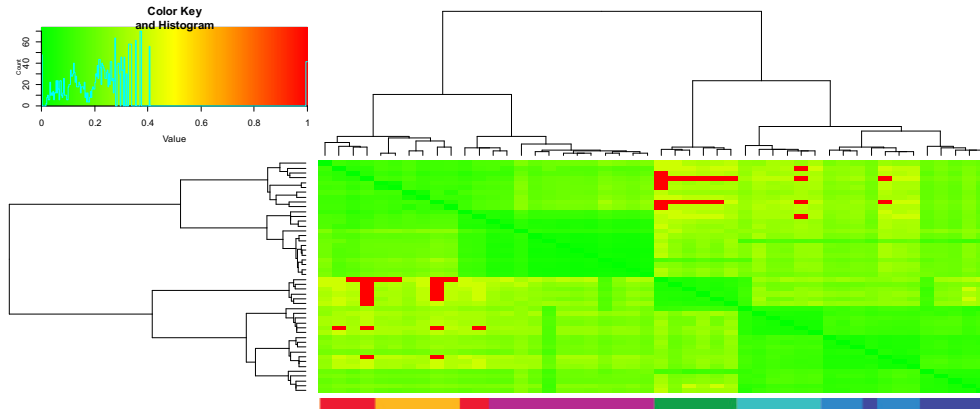
C- Simka, abundance Jaccard



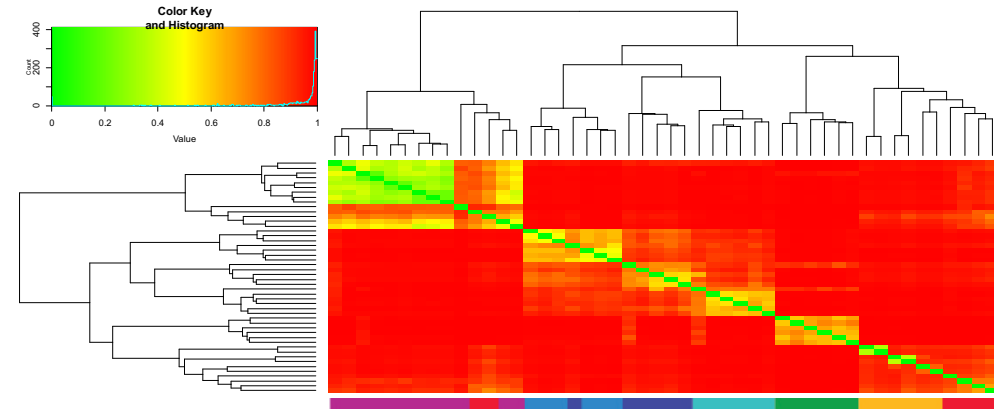
D- Simka, Abundance Bray-curtis



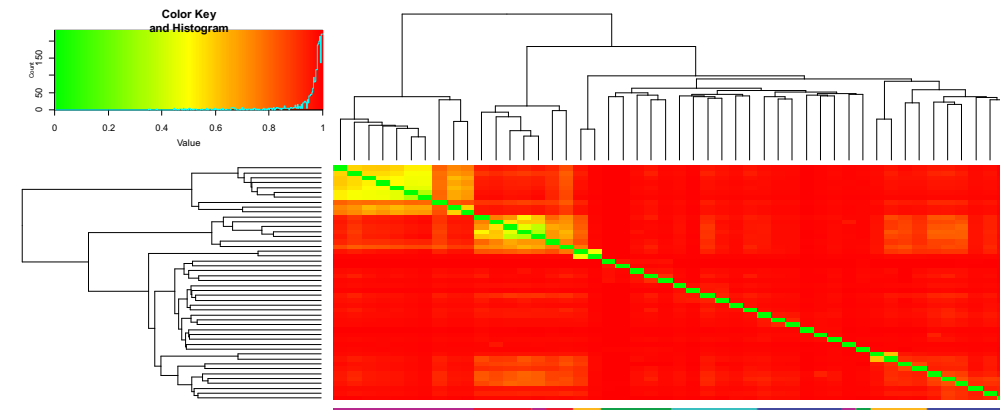
A - MASH



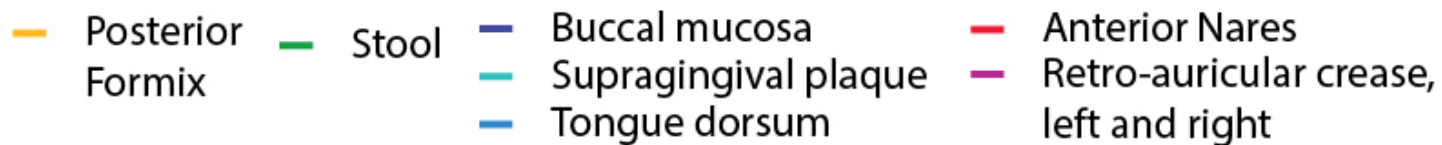
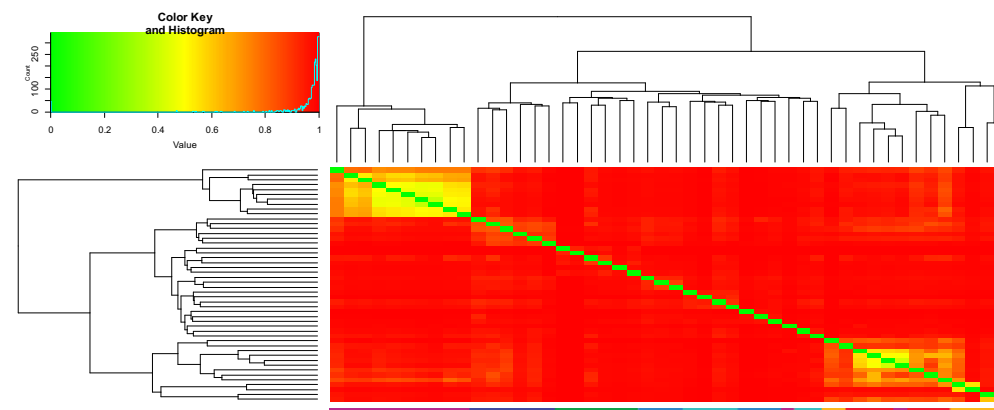
B - LIBRA, log weighting



C- Simka, abundance Jaccard



D- Simka, abundance Bray-curtis

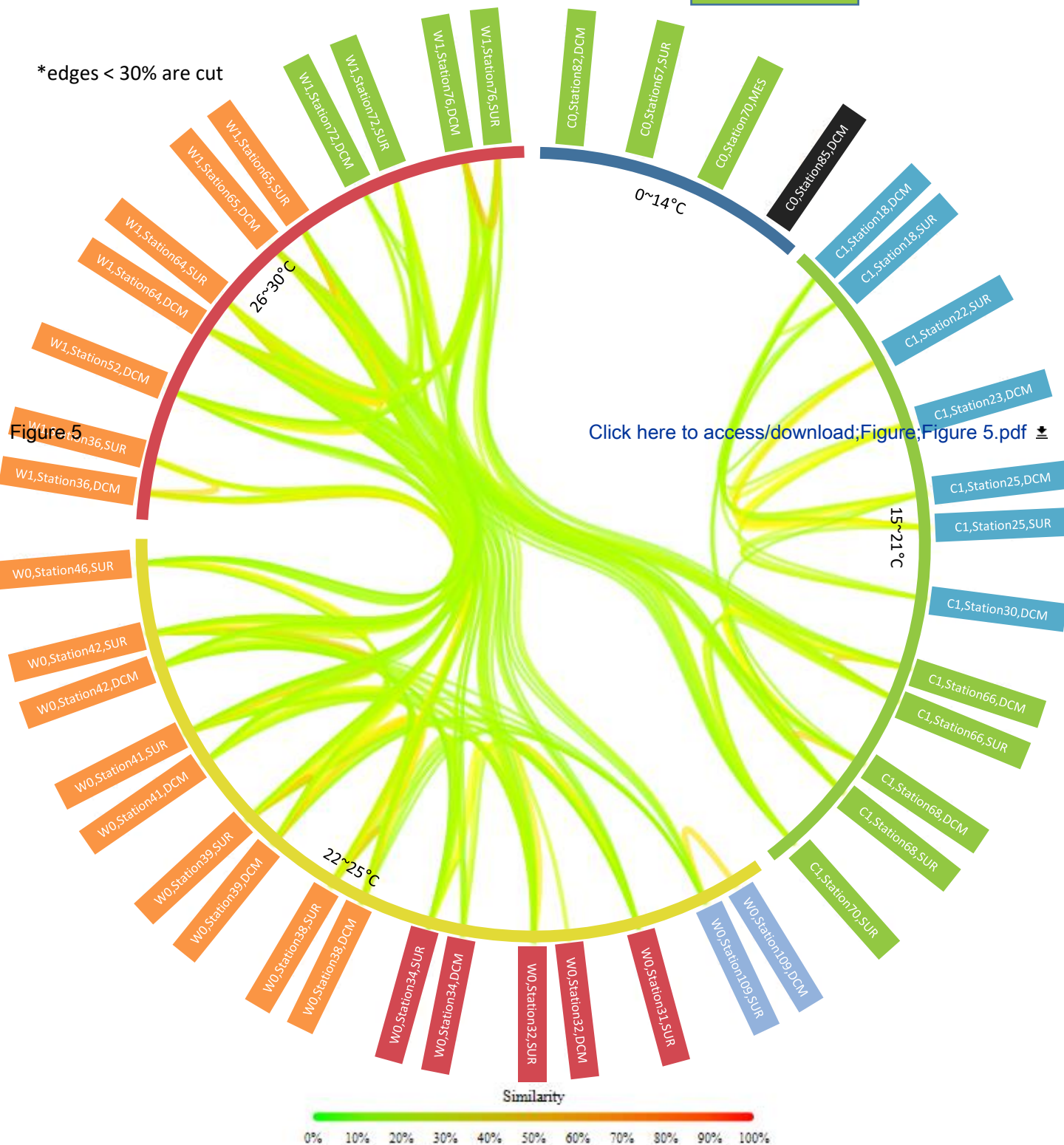


C0 - 0~14°C
 C1 - 15~21°C
 W0 - 22~25°C
 W1 - 26~30°C

SUR - surface
 DCM - deep chlorophyll maximum
 MES - mesopelagic

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

*edges < 30% are cut



Runtimes of Libra

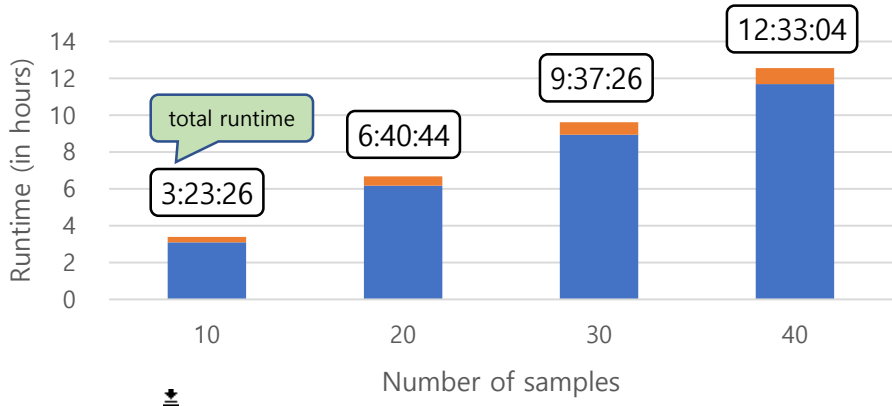
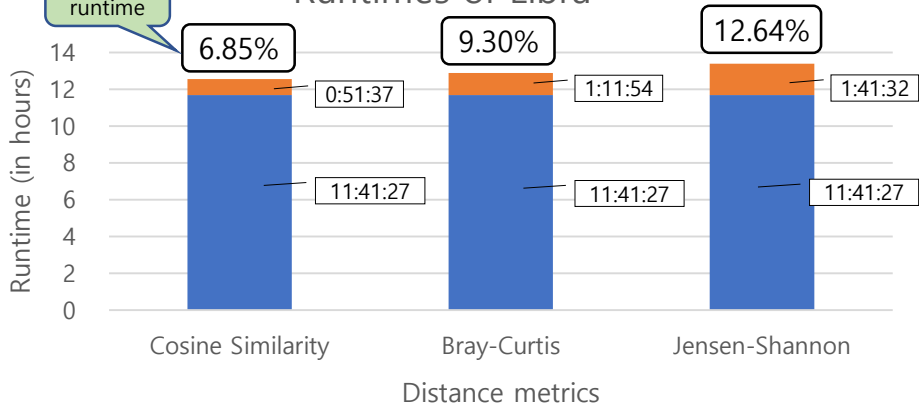


Figure 6.pdf

index construction

distance-matrix computation

Runtimes of Libra





Click here to access/download

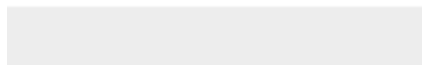
Supplementary Material

Supplemental Table1 -Comparable_tools.xlsx



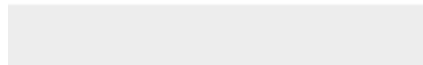


Click here to access/download
Supplementary Material
Supplemental Table 2.xlsx





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



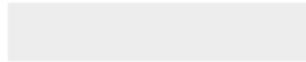


Click here to access/download
Supplementary Material
Supplemental Table 4.xlsx





Click here to access/download
Supplementary Material
Supplemental Table 5.xlsx





Click here to access/download
Supplementary Material
supplemental_Fig1.pdf





Click here to access/download
Supplementary Material
Supplemental Fig2.pdf



Click here to access/download
Supplementary Material
Supplemental Fig3.pdf





Click here to access/download
Supplementary Material
Supplemental Fig4.pdf





Click here to access/download
Supplementary Material
Supplemental Fig5.pdf

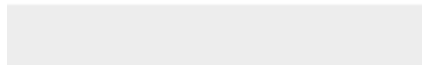




[Click here to access/download](#)

Supplementary Material

[Gigascience_Libra_Manuscript_TrackedChanges.docx](#)





Click here to access/download

Supplementary Material

Gigascience_Libra_Manuscript_TrackedChanges2.docx





Click here to access/download

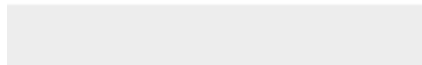
Supplementary Material

Gigascience_Libra_Manuscript_rev2_trackchanges.docx





Click here to access/download
Supplementary Material
Response_to_reviews_round2.docx





College of Agriculture
and Life Sciences

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,

A handwritten signature in black ink, appearing to read 'Bonnie Hurwitz', on a light-colored background.

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