

# GigaScience

## Optimized Distributed Systems Achieve Significant Performance Improvement on Sorted Merging of Massive Omics Data --Manuscript Draft--

|  |   |
|--|---|
| <b>Manuscript Number:</b>                            | GIGA-D-17-00267   |
| <b>Full Title:</b>                                   | Optimized Distributed Systems Achieve Significant Performance Improvement on Sorted Merging of Massive Omics Data   |
| <b>Article Type:</b>                                 | Technical Note  |
| <b>Funding Information:</b>                          |   |
| <b>Abstract:</b>                                     | <p>Background: Sorted merging of genomic data is a common data operation necessary in whole genome sequencing studies. It involves sorting and merging genomic data from different subjects by genomic locations. With the rapid increase of high throughput experimental data, the computational burden makes traditional methods designed for a single computer no longer feasible to this problem. The newly emerged distributed systems have the potential to offer a much needed boost in performance. However, carefully designed optimization schemas are required to take advantage of the increased computing power while overcoming bottlenecks to achieve maximum performance.</p> <p>Findings: In this study, we custom design optimized schemas for three Apache big data platforms, MapReduce, HBase and Spark, to perform sorted merging of massive genome-wide data. These schemas all adopt the divide-and-conquer strategy to split the merging job into sequential phases/stages consisting of subtasks which are conquered in an ordered, parallel and bottleneck-free way. In two illustrating examples, we test the performance of our schemas on merging multiple Variant Call Format (VCF) files into either a TPED or a VCF file, which are benchmarked with the traditional multiway-merge method and the popular VCFTools.</p> <p>Conclusions: Our experiments suggest that all three schemas deliver a significant performance improvement over existing methods. More importantly, they all show good scalability on input size and computing resources. Therefore our findings provide generalized scalable schemas for performing sorted merging on genetics and genomics data using these Apache distributed systems.</p> |
| <b>Corresponding Author:</b>                         | Zhaohui Qin<br><br>UNITED STATES  |
| <b>Corresponding Author Secondary Information:</b>   |   |
| <b>Corresponding Author's Institution:</b>           |   |
| <b>Corresponding Author's Secondary Institution:</b> |   |
| <b>First Author:</b>                                 | Xiaobo Sun  |
| <b>First Author Secondary Information:</b>           |   |
| <b>Order of Authors:</b>                             | Xiaobo Sun<br>Zhaohui Qin<br>Fusheng Wang<br>Jingjing Gao<br>Peng Jin   |
| <b>Order of Authors Secondary Information:</b>       |   |
| <b>Opposed Reviewers:</b>                            |   |

| <b>Additional Information:</b>  |                 |
|---|-----------------|
| <b>Question</b>   | <b>Response</b> |
| Are you submitting this manuscript to a special series or article collection?   | No              |
| <p><b>Experimental design and statistics</b></p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. 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             |
| <p><b>Resources</b></p> <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 <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>                     | Yes             |
| <p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (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 <a href="#">Minimum Standards Reporting Checklist</a>?</p> | Yes             |

1 **1 Optimized Distributed Systems Achieve Significant Performance**

2 **2 Improvement on Sorted Merging of Massive Omics Data**

3 **Xiaobo Sun<sup>1</sup>, Jingjing Gao<sup>2</sup>, Peng Jin<sup>3</sup>, Fusheng Wang<sup>4\*</sup>, Zhaohui Qin<sup>2,5\*</sup>**

4

5 <sup>1</sup>Department of Computer Sciences, Emory University, Atlanta, GA 30322, USA.

6 <sup>2</sup>Department of Medical Informatics, Emory University School of medicine, Atlanta, GA 30322, USA.

7 <sup>3</sup>Department of Human Genetics, Emory University School of Medicine, Atlanta, GA 30322, USA.

8 <sup>4</sup>Department of Biomedical Informatics, Stony Brook University, Stony Brook, NY 11794, USA.

9 <sup>5</sup>Department of Biostatistics, Emory University, Atlanta, GA 30322, USA.

10 X.S. Email: [xsun28@emory.edu](mailto:xsun28@emory.edu)

11 J.G. Email: [jingjing.gao@abbvie.com](mailto:jingjing.gao@abbvie.com)

12 P.J. Email: [peng.jin@emory.edu](mailto:peng.jin@emory.edu)

13 F.W. Email: [fusheng.wang@stonybrook.edu](mailto:fusheng.wang@stonybrook.edu)

14 Z.Q. Email: [zhaohui.qin@emory.edu](mailto:zhaohui.qin@emory.edu)

15 \*Correspondence: [zhaohui.qin@emory.edu](mailto:zhaohui.qin@emory.edu), [fusheng.wang@stonybrook.edu](mailto:fusheng.wang@stonybrook.edu)

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

16 **Abstract**

17 **Background:** Sorted merging of genomic data is a common data operation necessary in whole  
18 genome sequencing studies. It involves sorting and merging genomic data from different subjects  
19 by genomic locations. With the rapid increase of high throughput experimental data, the  
20 computational burden makes traditional methods designed for a single computer no longer feasible  
21 to this problem. The newly emerged distributed systems have the potential to offer a much needed  
22 boost in performance. However, carefully designed optimization schemas are required to take  
23 advantage of the increased computing power while overcoming bottlenecks to achieve maximum  
24 performance.

25 **Findings:** In this study, we custom design optimized schemas for three Apache big data platforms,  
26 MapReduce, HBase and Spark, to perform sorted merging of massive genome-wide data. These  
27 schemas all adopt the divide-and-conquer strategy to split the merging job into sequential  
28 phases/stages consisting of subtasks which are conquered in an ordered, parallel and bottleneck-  
29 free way. In two illustrating examples, we test the performance of our schemas on merging  
30 multiple Variant Call Format (VCF) files into either a TPED or a VCF file, which are  
31 benchmarked with the traditional multiway-merge method and the popular VCFTools.

1 32 **Conclusions:** Our experiments suggest that all three schemas deliver a significant performance  
2  
3  
4  
5 33 improvement over existing methods. More importantly, they all show good scalability on input  
6  
7  
8  
9 34 size and computing resources. Therefore our findings provide generalized scalable schemas for  
10  
11  
12  
13 35 performing sorted merging on genetics and genomics data using these Apache distributed systems.  
14  
15

16 36 **Keywords:** Sorted merging, whole genome sequencing, MapReduce, Hadoop, HBase, Spark.  
17  
18  
19  
20

21 37  
22  
23

## 24 38 **Findings**

### 25 39 **Introduction**

26  
27  
28  
29 40 With rapid development of high-throughput biotechnologies, genetics studies have entered the Big  
30  
31  
32  
33 41 Data era. Studies like Genome Wide Association Studies (GWASs), Whole Genome Sequencing  
34  
35  
36  
37 42 (WGS) and whole exome sequencing (WES) studies have produced a massive amount of data.  
38  
39  
40  
41  
42

43 43 The ability to efficiently process such massive data becomes increasingly important in a  
44  
45  
46  
47 44 successful large scale genetics study [1, 2]. Traditional single machine based methods are no  
48  
49  
50  
51  
52 45 longer feasible to process such big data due to the prohibitive computation time and I/O  
53  
54  
55  
56 46 bottleneck. It becomes increasingly attractive for investigators to take advantage of the powerful  
57  
58  
59  
60 47 distributed computing resources or the cloud to perform data processing and analyses [3]. Apache  
61  
62  
63  
64  
65

1 48 Foundation has been a leading force in this endeavor and has developed multiple platforms and  
2  
3  
4  
5 49 systems including Hadoop [4, 5], HBase [6] and Spark [7]. All these three Apache platforms have  
6  
7  
8  
9 50 gained increasing popularity in recent years, and have been endorsed and supported by major  
10  
11  
12  
13 51 vendors such as Amazon Web Services (AWS).  
14  
15  
16  
17 52  
18  
19  
20  
21 53 In bioinformatics, researchers have recently started to embrace distributed systems to process large  
22  
23  
24  
25 54 amount of high throughput omics data. For example, both the CloudBurst [8] and Crossbow  
26  
27  
28  
29 55 software [9] takes advantage of the Hadoop framework to accelerate sequencing read mapping and  
30  
31  
32  
33 56 SNP calling. The Collaborative Genomic Data Model (CGDM) [10] uses HBase to boost the  
34  
35  
36  
37 57 querying speed for the main classes of queries on genomic databases. The ADAM project [1],  
38  
39  
40  
41 58 built on the Spark platform, adapts the Sequence/Binary Alignment/Map (SAM/BAM) formats to  
42  
43  
44  
45 59 distributed computing environments. Industry cloud computing vendors such as Amazon [11] and  
46  
47  
48  
49 60 Google [12] are also beginning to provide specialized environments to ease genomics data  
50  
51  
52  
53 61 processing in the cloud.  
54  
55  
56  
57 62  
58  
59  
60  
61  
62  
63  
64  
65

1 63 Despite their potentials, applications of Apache big data platform in genetics and genomics studies  
2  
3  
4  
5 64 are still relatively limited. We believe there are plenty of opportunities as data becomes larger and  
6  
7  
8  
9 65 more complex. One particular example is sorted merging, which is a ubiquitous operation in  
10  
11  
12  
13 66 processing genetics and genomics data. As an example, in WGS, variants identified from  
14  
15  
16  
17 67 individuals are often called and stored in separate VCF files, subsequently these VCF files need to  
18  
19  
20  
21 68 be merged (into a VCF or TPED file) as required by downstream analyses such as PLINK [13]  
22  
23  
24  
25 69 and BlueSNP [14, 15]. Either a VCF or TPED file requires data to be sorted by genomic location,  
26  
27  
28  
29 70 thus these tasks are equivalent to the well-known sorted full-outer-joining problem [16, 17].  
30  
31  
32  
33 71 Currently, they are handled by software such as VCFTools [18] and PLINK. These utilities  
34  
35  
36  
37 72 become very cumbersome even in the face of a moderate scale of genomic data. The main reason  
38  
39  
40  
41 73 is that most of these tools adopt the multiway-merge-like method [19] with a priority queue as the  
42  
43  
44  
45 74 underlying data structure to ensure the output order. A key deficiency of such method is that it can  
46  
47  
48  
49 75 only have one consumer to access items from the queue, which literally makes it single-threaded,  
50  
51  
52  
53 76 even if there can be parallel producers that put items into the queue. Therefore, these single-  
54  
55  
56  
57 77 machine based tools are inefficient and time-consuming when handling large datasets.  
58  
59  
60  
61  
62  
63  
64  
65

1 78  
2  
3  
4  
5  
6 79 In this study, we use the case of the sorted-merging of multiple VCF files to a single file to  
7  
8  
9  
10 80 demonstrate the benefits of using distributed platforms. However, simply running sorted merging  
11  
12  
13  
14 81 on a distributed system runs into problems of bottlenecks, hotspots and unordered results  
15  
16  
17  
18 82 commonly seen in parallel computations. Rather, we believe working schemas custom designed  
19  
20  
21  
22 83 for each specific distributed platform are required to unleash the full potential of these distributed  
23  
24  
25  
26 84 systems. We propose and implement three schemas running on Hadoop, Spark and HBase  
27  
28  
29  
30 85 respectively to overcome the limitations of both single-machine and simple distributed system  
31  
32  
33  
34 86 based methods. We choose these three platforms because they are representative cloud distributed  
35  
36  
37  
38 87 systems providing data partitioning based parallelism with distributed storage, data partitioning  
39  
40  
41  
42 88 based parallelism with in-memory based processing, and high dimensional table like distributed  
43  
44  
45  
46 89 storage, respectively. Hadoop [4] is the open source implementation of MapReduce [5] based  
47  
48  
49  
50 90 parallel key-value processing technique, and has the advantage of transparency and simplicity.  
51  
52  
53  
54 91 HBase [6] is a data warehousing platform which adopts Google's BigTable data storing structure  
55  
56  
57  
58 92 [20] to achieve high efficiency in storing and reading/writing large scale of sparse data. Spark [7]  
59  
60  
61  
62  
63  
64  
65

1 93 introduces the concept of Resilient Distributed Dataset (RDD) and Directed Acyclic Graph (DAG)  
2  
3  
4  
5 94 execution to parallel key-value processing, thus enabling fast, robust and repetitive in-memory  
6  
7  
8  
9 95 data manipulations. Specifically, our schemas involve dividing the job into multiple phases  
10  
11  
12  
13 96 corresponding to tasks of loading, mapping, filtering, sampling, partitioning, shuffling, merging  
14  
15  
16  
17 97 and outputting. Within each phase, data and tasks are evenly distributed across the cluster,  
18  
19  
20  
21 98 enabling processing large scale of data in a parallel and scalable manner, which in turn  
22  
23  
24  
25 99 significantly boosts performance.  
26  
27  
28

29 100  
30  
31

## 32 101 **Methods**

### 33 102 **Overview**

34  
35  
36  
37  
38  
39  
40  
41 103 Compared to using the multiway-merge method [19] or a relational database based approach, the  
42  
43  
44  
45 104 benefits of using the three Apache distributed platforms to perform sorted merging are three-fold.  
46  
47  
48  
49 105 First, with representation of genomic locations as keys and genotypes as values, it is readily  
50  
51  
52  
53 106 transformed into the key-value model on which all three platforms offer a rich set of parallel  
54  
55  
56  
57 107 operations. Second, data in VCF files are semi-structured. Semi-structured data ideally fit for all  
58  
59  
60  
61  
62  
63  
64  
65

1 108 three platforms which allow defining the schema during data loading, avoiding the preprocessing  
2  
3  
4  
5 109 of raw data into a rigid schema as in a relational database. Third, the merged results are outputted  
6  
7  
8  
9 110 onto a distributed file system such as HDFS and Amazon S3 which can be directly used for  
10  
11  
12  
13 111 subsequent cluster-based GWAS or WGS analytical tools such as BlueSNP.  
14  
15  
16  
17 112  
18  
19  
20  
21 113 Despite these advantages, simply performing sorted merging on distributed systems will not  
22  
23  
24  
25 114 deliver expected results for the following reasons. First, it can lead to globally unsorted results.  
26  
27  
28  
29 115 Hash-based shuffling of input data is the default mechanism for distributing data to parallel  
30  
31  
32  
33 116 working units in the system. However, shuffling will lead to globally unsorted results. Second,  
34  
35  
36  
37 117 bottleneck and hotspot can happen during the processing in the cluster. Bypassing the hashing  
38  
39  
40  
41 118 based shuffling can lead to unbalanced workload across the cluster, result in straggling computing  
42  
43  
44  
45 119 units which become the bottlenecks for response time. In addition, for parallel loading of presorted  
46  
47  
48  
49 120 data into HBase, data being loaded from all the loading tasks will hit the same node  
50  
51  
52  
53 121 simultaneously while other machines are idling, leading to an I/O hotspot. Third, sampling costs  
54  
55  
56  
57 122 could become prohibitive. Although Hadoop provides a native utility named *total-order-merging*  
58  
59  
60 123 [16] to achieve both workload balance and global order, it involves transferring to and sampling  
61  
62  
63  
64  
65

1 124 all the data onto a single node. The communication cost over the network and disk I/O can be  
2  
3  
4  
5 125 prohibitive when data size is very large. In the following sections, we will illustrate how our  
6  
7  
8  
9 126 custom designed schema are able to overcome these limitations in detail.  
10  
11  
12  
13 127  
14  
15  
16  
17 128 **Data Formats and Operations**  
18  
19  
20  
21 129 In a typical WGS, data analysis often starts from individual genotype files in VCF format [21]. A  
22  
23  
24  
25 130 VCF file contains data arranged into a table consisting of eight mandatory fields including  
26  
27  
28  
29 131 chromosome (CHROM), the genomic coordinate of the start of the variant (POS), the reference  
30  
31  
32  
33 132 allele (REF), a comma separated list of alternate alleles (ALT), among others. In our experiments,  
34  
35  
36  
37 133 we use a dataset consisting of the VCF files of 93 individuals [22] generated from Illumina's  
38  
39  
40  
41 134 BaseSpace software (Left tables in Figure 1). Each file has around 4-5 million rows, each  
42  
43  
44  
45 135 representing one of the individual's genomic variants, with a size of about 300 megabytes. In an  
46  
47  
48  
49 136 attempt to protect the privacy of the study subjects, we apply the following strategy to conceal  
50  
51  
52  
53 137 their real genetic variant information contained in the VCF files: we first transform each original  
54  
55  
56  
57 138 genomic location by multiplying it with an undisclosed constant real number, taking the floor  
58  
59  
60  
61 139 integer of the result, and then add another undisclosed constant integer number.  
62  
63  
64  
65

1 140  
2  
3  
4  
5 141 It is common that multiple VCF files need to be merged into a single TPED file for analysis tools  
6  
7  
8  
9 142 such as PLINK. A TPED file resembles a big table, aggregating genotypes of all individuals under  
10  
11  
12  
13 143 investigation by genomic location (Right table in Figure 1). The merging follows several rules.  
14  
15  
16  
17 144 First, records having an unqualified filter value are discarded. Second, genotypes in VCF files are  
18  
19  
20  
21 145 stored as binary codes where 0 stands for reference allele while 1 stands for mutant allele. Binary  
22  
23  
24  
25 146 codes must be translated into corresponding types of nucleotides in the TPED file. Third, all  
26  
27  
28  
29 147 individuals need to have a genotype for genomic locations that appears in at least one VCF file.  
30  
31  
32  
33 148 The default genotype for missing values is homozygous reference alleles.  
34  
35

36 149  
37  
38  
39

## 40 150 **MapReduce Schema**

41  
42  
43  
44 151 MapReduce [5] is a parallel computing model based on a *split-apply-combine* strategy for data  
45  
46  
47  
48 152 analysis, in which data are mapped to key-values for splitting (mapping), shuffling and combining  
49  
50  
51  
52 153 (reducing) for final results. We use Apache Hadoop-2.7 as the system for our implementation. Our  
53  
54  
55  
56 154 optimized schema consists of two MapReduce phases, as shown in Figure 2.  
57  
58  
59

60 155  
61  
62  
63  
64  
65

1 156 *First MapReduce phase.* Raw data are loaded from HDFS into parallel mappers to perform the  
2  
3  
4  
5 157 following tasks: First, unqualified data are filtered out and qualified ones are mapped to key-value  
6  
7  
8  
9 158 pairs. The mapper output key is a genomic location and output value is genotype and individual  
10  
11  
12  
13 159 ID. Second, Key-value pairs are grouped together by their chromosome and temporarily saved as  
14  
15  
16  
17 160 compressed Hadoop sequence files [23] for faster I/O in the second MapReduce phase. With this  
18  
19  
20  
21 161 grouping, we can merge records from selected chromosomes of interests rather than from all of  
22  
23  
24  
25 162 them. Meantime, these records are sampled to explore their key distribution profile along the  
26  
27  
28  
29 163 chromosomes for determining boundaries in between each pair of which there is approximately an  
30  
31  
32  
33 164 equal number of records. Specifically, the genomic locations of sampled-out records for each  
34  
35  
36  
37 165 chromosome are used as boundaries to split the chromosome into disjoint segments. Because  
38  
39  
40  
41 166 records falling in the same segment will be assigned to the same reducer in the later phase,  
42  
43  
44  
45 167 boundaries calculated in this way ensure that the workload of each reducer is balanced. There are  
46  
47  
48  
49 168 two rounds of samplings. The first one happens in each mapper with a pre-specified sampling rate,  
50  
51  
52  
53 169 which in our case is set to 0.0001. To separate sampled records by chromosome they are  
54  
55  
56  
57 170 distributed to different reducers in this phase based on their chromosomes, where they are sampled  
58  
59  
60  
61 171 again with a rate equal to the reciprocal of input file number. This second sampling limits the  
62  
63  
64  
65

1 172 number of final sampled records even in the face of a large number of input files. Because the  
2  
3  
4  
5 173 number of reducers instantiated in the second phase is decided by the number of sampled records,  
6  
7  
8  
9 174 we can therefore avoid launching unnecessary reducers thus reducing task overhead.  
10  
11  
12  
13 175  
14  
15  
16  
17 176 *Second MapReduce phase.* In this phase, multiple parallel MapReduce jobs are created, and each  
18  
19  
20  
21 177 job specifically handles all records of a single chromosome outputted as sequence files in the first  
22  
23  
24  
25 178 phase. Within each job, a partitioner shuffles records to the appropriate reducer by referring to the  
26  
27  
28  
29 179 boundaries from the previous phase, so that records falling in between the same pair of boundaries  
30  
31  
32  
33 180 are aggregated together. Finally, each reducer sorts and merges aggregated records by genomic  
34  
35  
36  
37 181 location before outputting them to a TPED file. In this way, globally sorted merging can be  
38  
39  
40  
41 182 fulfilled.  
42  
43  
44  
45 183  
46  
47  
48 184 **HBase Schema**  
49  
50  
51  
52 185 HBase [6] is a column-oriented database where data are grouped into column families and split  
53  
54  
55  
56 186 horizontally into regions spreading across the cluster. With this data storing structure, it supports  
57  
58  
59  
60 187 efficient sequential reading and writing of large-scale data as well as fast random data accessing.  
61  
62  
63  
64  
65

1 188 Also, HBase is storage efficient because it can remember null values without saving them on disk.  
2  
3  
4  
5 189 These features make HBase an ideal platform for managing large, sparse data with relatively low  
6  
7  
8  
9 190 latency which naturally fits the sorted merging case. We use the HBase-1.3 as the system for our  
10  
11  
12  
13 191 implementation. As shown in Figure 3, our optimized HBase schema is divided into three phases  
14  
15  
16  
17 192 as discussed next.  
18  
19  
20  
21 193  
22  
23  
24  
25 194 **1) Sampling phase**  
26  
27  
28  
29 195 The main challenge of HBase lies in that it is not uncommon to find that one server of the cluster  
30  
31  
32  
33 196 becomes a computational hotspot. This can happen when it starts loading a table from a single  
34  
35  
36  
37 197 region hosted by a single node. Therefore, we need to presplit the table into regions of  
38  
39  
40  
41 198 approximately equal size before loading. The sampling phase is introduced to determine  
42  
43  
44  
45 199 reasonable presplitting regional boundaries. The total region number is set to be half of the  
46  
47  
48  
49 200 number of input files so that the size of each region is approximately 1GB. Meanwhile, mappers  
50  
51  
52  
53 201 of this phase also output qualified records as compressed Hadoop sequence files on HDFS which  
54  
55  
56  
57 202 are used as inputs in the next phase. In addition, filtering and key-value mapping also take place in  
58  
59  
60  
61 203 this phase.  
62  
63  
64  
65

1 204  
2  
3  
4  
5 205 2) Bulk loading phase  
6  
7  
8  
9 206 Even when the table has been presplit evenly, the hotspot problem of loading sorted inputs is not  
10  
11  
12  
13 207 yet fully solved because sorted records are loaded sequentially and all the records being loaded at  
14  
15  
16  
17 208 any instant still hit the same region and server. This necessitates the adding of this phase. During  
18  
19  
20  
21 209 the bulk loading, the key and value of each record outputted from previous phase is converted into  
22  
23  
24  
25 210 HBase's binary row-key and column-value respectively, and saved as HFile, HBase's native  
26  
27  
28  
29 211 storage format. The row-key here is in the form of chromosome-genomic location, and column-  
30  
31  
32  
33 212 value refers to reference allele, individual ID and genotype. The bulk loading populates each  
34  
35  
36  
37 213 HFile with records falling in the same pair of presplit regional boundaries. Because HFiles are  
38  
39  
40  
41 214 written simultaneously by parallel mappers/reducers, all working nodes are actively involved and  
42  
43  
44  
45 215 the regional hotspot is thus circumvented. Upon finishing writings, the HBase can readily load  
46  
47  
48  
49 216 HFiles in parallel into the table by simply moving them into local storage folders. This procedure  
50  
51  
52  
53 217 is therefore at least a magnitude faster than the normal loading. The order of records in the table is  
54  
55  
56  
57 218 guaranteed because they are internally sorted by writing reducers and HBase's Log-Structured  
58  
59  
60  
61 219 Merge-tree [24]. It is noteworthy to mention that VCF records are sparse, thus HBase is very  
62  
63  
64  
65

1 220 storage-efficient.  
2  
3  
4  
5 221  
6  
7  
8  
9 222 **3) Exporting Phase**  
10  
11  
12  
13 223 A scan is performed on the table. It involves launching parallel mappers each receiving records  
14  
15  
16  
17 224 from a HBase region, filling in missing genotypes, concatenating records with the same row-key,  
18  
19  
20  
21 225 and outputting final results into TPED files.  
22  
23  
24  
25 226  
26  
27  
28  
29 227 **Spark Schema**  
30  
31  
32  
33 228 Spark [7] is a distributed engine that embraces the ideas of MapReduce and Resilient Distributed  
34  
35  
36  
37 229 Dataset (RDD). It can save intermediate results in the form of RDD in memory, and perform  
38  
39  
40  
41 230 computation on them. Also, its computations are lazily evaluated, which means the execution plan  
42  
43  
44  
45 231 can be optimized since it tries to include as many computational steps as possible. As a result, it is  
46  
47  
48  
49 232 ideal for iterative computations such as sorted merging. We implement our optimized Spark  
50  
51  
52  
53 233 schema on Spark-2.1. It has three stages as shown in Figure 4. Stage I involves loading raw data  
54  
55  
56  
57 234 as RDDs, filtering, and mapping RDDs to paired-RDDs with keys (chromosome-genomic  
58  
59  
60  
61 235 position) and values (reference allele, individual ID and genotype). This stage ends with a sort-by-  
62  
63  
64  
65

1 236 key shuffling that repartitions and sorts PairRDD records so that records with the same key are  
2  
3  
4  
5 237 aggregated together. In Stage II, aggregated PairRDD records of the same key are merged into  
6  
7  
8  
9 238 TPED format and converted back to RDD records for outputting. However, Spark's native family  
10  
11  
12  
13 239 of group-by-key functions cannot be used here because their default partitioner is hash-based and  
14  
15  
16  
17 240 different from the range-based partitioner used by previous sort-by-key function. Consequently,  
18  
19  
20  
21 241 the merged results would be reshuffled into an unsorted status. We therefore optimize the merging  
22  
23  
24  
25 242 to bypass these functions, being performed locally without data reshuffling to ensure both order  
26  
27  
28  
29 243 and high speed. Finally in Stage III, merged RDD records are saved as TPED files.  
30  
31  
32  
33 244  
34  
35  
36  
37 245 Execution parallelism has an important impact on the performance. To maximize performance, the  
38  
39  
40  
41 246 number of parallel tasks is set to the number of input files. In this way, the data locality is  
42  
43  
44  
45 247 maximized and each task is assigned a proper amount of work. In addition, unlike using  
46  
47  
48  
49 248 MapReduce or HBase, when performing sorting by key, no explicit sampling is needed because  
50  
51  
52  
53 249 Spark keeps track of the number of records before determining repartition boundaries.  
54  
55  
56  
57 250  
58  
59

## 60 251 **Results**

1 252 We conduct all experiments using Amazon’s Elastic MapReduce (EMR) service. Within the  
2  
3  
4  
5 253 infrastructure, we choose EC2 working nodes of m3.xlarge type, which has four High Frequency  
6  
7  
8  
9 254 Intel Xeon E5-2670 v2 (Ivy Bridge) Processors and 15GB memory. We use a dataset consisting of  
10  
11  
12  
13 255 the VCF files of 93 individuals [22] generated from Illumina's BaseSpace software.  
14  
15

16  
17 256

### 21 257 **Overall Performance Analysis of Clustered-based Schemas**

22  
23  
24  
25 258 Our primary goal is to explore the scalability of the three schemas on input data size and available  
26  
27  
28  
29 259 computing resources, namely CPUs. To achieve this, in this experiment we adjust the number of  
30  
31  
32  
33 260 input files from 10 to 93, with an approximate total uncompressed size from 2.5 G to 20 G, and  
34  
35  
36  
37 261 conduct the experiment using a varying number of working nodes from 3 to 18, namely 12 to 72  
38  
39  
40  
41 262 cores.  
42  
43

44 263

45  
46  
47  
48 264 As Figure 5 shows, for all three schemas, given a fixed number of cores, the execution time  
49  
50  
51  
52 265 increases linearly with the increased number of input files. On the one hand, the increasing trend is  
53  
54  
55  
56 266 apparent with fewer cores because each core is fully utilized and the more input files, the larger  
57  
58  
59  
60 267 number of parallel tasks are assigned to it. For example, given 12 cores, as the file number  
61  
62  
63  
64  
65

1 268 increases from 10 to 93, the execution time increases from 739 to 2,281 seconds for the  
2  
3  
4  
5 269 MapReduce schema, from 375 to 2,751 seconds for the HBase schema, and from 361 to 1,699  
6  
7  
8  
9 270 seconds for the Spark schema, respectively. On the other hand, with relatively more cores such as  
10  
11  
12  
13 271 72, this linear increasing trend is less pronounced because there are more cores than tasks so that  
14  
15  
16  
17 272 all cores are assigned at most one task. We also notice that when input file size is small to  
18  
19  
20  
21 273 moderate, the Spark schema does not always show consistent improvement in terms of execution  
22  
23  
24  
25 274 time when using more cores, for example, the intersection of curves of 24 and 72 cores in Figure  
26  
27  
28  
29 275 5c. This phenomenon is attributed to the limitation of Spark's internal task assignment policy  
30  
31  
32  
33 276 which gives rise to the possibility that some nodes are assigned more than one tasks while others  
34  
35  
36  
37 277 remain idle.  
38  
39  
40  
41 278  
42  
43  
44  
45 279 In another experiment in which the input file number is fixed at 93, the core number increases  
46  
47  
48  
49 280 from 12 to 72 (Figure 6). For all three schemas, execution time is reduced with more cores, from  
50  
51  
52  
53 281 2,281 to 514 seconds for MapReduce, from 2,751 to 591 seconds for HBase, and from 1,699 to  
54  
55  
56  
57 282 460 seconds for Spark, respectively. Therefore, all three schemas demonstrate nice scalability on  
58  
59  
60  
61 283 input data size and computing resources.  
62  
63  
64  
65

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

284

## 285 **The Anatomic Performances Analysis of Cluster-based Schemas**

286 Another important goal of our study is to identify potential performance bottlenecks, so we  
287 evaluate the execution time of each phase/stage of all three schemas. Figure 7 shows the trend of  
288 anatomic computing time spent on merging increasing number of VCF files using 48 cores. For  
289 the MapReduce schema (Figure 7a), its two phases account for a comparable proportion of total  
290 time and both show a linear or sublinear increasing pattern. For the three phases of the HBase  
291 schema (Figure 7b), they generally scale well with the input file number. Meanwhile, the second  
292 phase becomes more dominant with more input files owing to the larger amount of shuffled data  
293 during the writing of HFiles. However, we do not consider it as a bottleneck since all tasks of this  
294 phase are parallelized with no workload or computational hotspot. We do not observe an obvious  
295 super-linear increment pattern from the figure either. Finally, Figure 7c shows the time costs of  
296 three stages of the Spark schema. They show a uniform increasing trend with input file number.  
297 Among them, the second one takes up a considerable proportion of the total execution time as it  
298 has the relatively expensive sort-by-key shuffling operation. Although no data is shuffled in the  
299 first stage, its time lapse is close to that of the second stage. This is because at the end of the first

1 300 stage, data are sampled for determining the boundaries used by sort-by-key's range partitioner.  
2  
3

4  
5 301 This operation demands a considerable execution time because it scans all the data and balances  
6

7  
8  
9 302 them if necessary.  
10

11  
12  
13 303

14  
15  
16  
17 304 Put together, these results suggest that all phases/stages of the three schemas scale well on input  
18

19  
20  
21 305 data size. Therefore we are not expecting to see any bottleneck when dealing with even larger  
22

23  
24  
25 306 scale of data.  
26

27  
28  
29 307  
30

### 31 308 **Comparisons between Single Machine Based Methods and Cluster-based Schemas** 32 33

34  
35  
36 309 Another intriguing question in our experiments is how much can our schemas outperform the  
37

38  
39  
40 310 current single machine based methods and applications. To achieve this, in the first experiment,  
41

42  
43  
44 311 we test the performance of merging 40 VCF files into one VCF using VCFTools v4.2 as a  
45

46  
47  
48 312 benchmark. As shown in Table 2, VCFTools takes 30,189 seconds while the fastest cluster-based  
49

50  
51  
52 313 schema, MapReduce-based, takes only 484 seconds using 72 cores, which is 62-fold faster. In the  
53

54  
55  
56 314 second experiment, we test the performance of merging of VCF files into a TPED using multiway-  
57

58  
59  
60 315 merge based implementation as a benchmark. We choose to implement multiway-merge because,  
61  
62  
63  
64  
65

1 316 to our best knowledge, currently no software/application is available to perform this task with a  
2  
3  
4  
5 317 simple call. Rather, VCF files need to be manually converted to individual TPED files first which  
6  
7  
8  
9 318 in turn are merged together using merging utility of PLINK which essentially is based on  
10  
11  
12  
13 319 multiway-merge. Our multi-way merge implementation saves this manual conversion step, and is  
14  
15  
16  
17 320 also more concise and efficient than PLINK. The multiway-merge implementation is tested on a  
18  
19  
20  
21 321 single node while the three schemas on a cluster of nodes with 72 cores. Initially with fewer input  
22  
23  
24  
25 322 files, the execution time difference is 399 seconds or about 2.8-fold between multiway-merge and  
26  
27  
28  
29 323 the fastest cluster-based schema, MapReduce. However, this difference becomes significant with  
30  
31  
32  
33 324 more input files. For example, the largest difference is 5,585 seconds or about 13-fold (Figure 8)  
34  
35  
36  
37 325 on merging 93 files. As an extreme test of merging 642 VCF files (not shown), the computing time  
38  
39  
40  
41 326 is 1,228 minutes for multiway merge implementation versus 11.3 minutes of the MapReduce  
42  
43  
44  
45 327 schema running on a 400-core cluster, more than 100-fold of speed up.  
46  
47  
48  
49 328  
50  
51  
52 329 We also compare the performances among the three schemas all of which are evaluated on a 72-  
53  
54  
55  
56 330 core cluster with increasing number of files as inputs (Figure 8). It turns out that the three schemas  
57  
58  
59  
60 331 have comparable performance. More specifically, MapReduce-based schema performs best with  
61  
62  
63  
64  
65

1 332 small input size, HBase-based schema performs best with moderate input size, while Spark-based  
2  
3  
4  
5 333 schema performs best with large input size. The rationale behind such an observation is when the  
6  
7  
8  
9 334 input data size is small, MapReduce can make the most usage of computing resources because it  
10  
11  
12  
13 335 has a constant 25 parallel jobs (one for each of chromosomes 1-22, X Y and M (Mitochondria)) in  
14  
15  
16  
17 336 its second phase. In contrast, Spark has much fewer tasks with a number equals to the number of  
18  
19  
20  
21 337 input files for achieving maximized data-task locality. When the input data size is moderate,  
22  
23  
24  
25 338 HBase triumphs due to its internal sorting and relative compact storage format of intermediate  
26  
27  
28  
29 339 data. When the input data size is large, Spark-based schema outperforms the other two owing to its  
30  
31  
32  
33 340 least number of data shuffling (only one), execution plan optimization, and ability to cache  
34  
35  
36  
37 341 intermediate results in memory. We caution that the computing time may fluctuate depending on  
38  
39  
40  
41 342 the genomic location profile of input files as well as the data loading balance of the HDFS.  
42  
43  
44  
45 343  
46  
47

## 344 **Discussion**

50  
51  
52 345 In this report, we describe three cluster-based schemas running on Apache Hadoop (MapReduce),  
53  
54  
55  
56 346 HBase and Spark platforms respectively for performing sorted merging of variants identified from  
57  
58  
59  
60 347 WGS. We manage to show that all three schemas are highly scalable on both input data size and  
61  
62  
63  
64  
65

1 348 computing resources, suggesting that large scale sequencing of variant data can be merged  
2  
3  
4  
5 349 efficiently given computing resources that are readily available in the cloud. We also show that  
6  
7  
8  
9 350 even with a moderate-sized cluster and input data, all three schemas are able to significantly  
10  
11  
12  
13 351 outperform the broadly-used, single-machine based VCFTools and multiway-merge  
14  
15  
16  
17 352 implementation. We expect a much more significant performance improvement when merging a  
18  
19  
20  
21 353 much larger scale of data using a larger cluster or the cloud.  
22  
23  
24  
25 354  
26  
27  
28  
29  
30 355 Unlike normal merging, efficient sorted merging of many large tables has always been a difficult  
31  
32  
33  
34 356 problem in the field of data management. Multiway-merge is the most efficient single-machine  
35  
36  
37  
38 357 based method for sorted merging, but its performance is limited by the disk I/O [25]. Sorted  
39  
40  
41  
42 358 merging also places challenges to distributed system based solutions because neither the efficient  
43  
44  
45  
46 359 hash-based merging nor caching the intermediate table in shared memory is feasible [26].  
47  
48  
49  
50 360 Although a utility named total-order-joining is provided by the Hadoop for addressing this  
51  
52  
53  
54 361 problem, it suffers from both network communication and local disk I/O bottleneck, thus is not  
55  
56  
57  
58 362 scalable [16, 27]. In contrast, our schemas divide this problem into different phases of tasks each  
59  
60  
61  
62  
63  
64  
65

1 363 conquered in parallel to bypass these bottlenecks and achieve maximum parallelism and  
2  
3  
4  
5 364 scalability. Furthermore, in addition to merging sequencing variant data, the schemas can be  
6  
7  
8  
9 365 generalized for other key-based, sorted merging problems that are frequently encountered in  
10  
11  
12  
13 366 genetics and genomics data processing. As an example, they can be slightly modified to merge  
14  
15  
16  
17 367 multiple BED format files such as ChIP-seq peak lists [28] and other genomic regions of interest.  
18  
19  
20  
21 368 Another potentially useful feature is that, unlike traditional sorted merging algorithms which  
22  
23  
24  
25 369 usually require presorted inputs for better performance, our schemas are free of such a  
26  
27  
28  
29 370 requirement.  
30  
31  
32  
33 371  
34  
35  
36  
37 372 Finally, in light of the different features and specialties of the three platforms, each of the three  
38  
39  
40  
41 373 schemas we developed has its own advantages and disadvantages in different application scenarios  
42  
43  
44  
45 374 as summarized in Table 1. For example, the MapReduce schema is good for static one-time, non-  
46  
47  
48  
49 375 incremental merging on small to moderate-sized data since it can have the most parallel jobs, the  
50  
51  
52  
53 376 least overhead, and the most transparent workflow. The HBase schema, supported by data  
54  
55  
56  
57 377 warehousing technologies, fits for incremental merging since it does not need to re-merge existing  
58  
59  
60  
61 378 results with new ones from the scratch only if the incremental merging is performed on the same  
62  
63  
64  
65

1 379 chromosomes. Also, it provides highly-efficient storage and On-Line Analytical Processing  
2  
3  
4  
5 380 (OLAP) on merged results. The Spark schema is ideal for merging large scale of data because it  
6  
7  
8  
9 381 has the least data shuffling and keeps intermediate results in memory. A bonus brought by Spark is  
10  
11  
12  
13 382 that subsequent statistical analyses can be carried out directly on the merged results using its rich  
14  
15  
16  
17 383 set of parallel statistical utilities.  
18  
19  
20

21 384

## 25 385 **Availability and Requirements**

27  
28  
29 386 Project name: CloudMerge  
30  
31

32  
33 387 Project home page: <https://github.com/xsun28/CloudMerge>  
34  
35

36  
37 388 Operating system(s): Linux  
38  
39

40  
41 389 Programming language: Java  
42  
43

44  
45 390 Other requirements: Java 1.7 or higher, Hadoop-2.7, HBase-1.3, Spark-2.1  
46  
47

48  
49 391 License: Apache License 2.0  
50  
51

52  
53 392

## 56 393 **Availability of Data and Materials**

57  
58  
59  
60  
61  
62  
63  
64  
65

1 394 The source code of the project is available on GitHub. The original 93 individual VCF files used  
2  
3  
4  
5 395 in our experiments are from the Consortium on Asthma among African-ancestry Population in the  
6  
7  
8  
9 396 Americas (CAAPA) [22]. To conceal the potential individual identifiable genotype information  
10  
11  
12  
13 397 from the public, we encrypt the authentic genomic location of all VCF files to generate a new  
14  
15  
16  
17 398 batch of encrypted VCF files, which are available on AWS S3 as  
18  
19  
20  
21 399 <https://s3.amazonaws.com/xsun316/encrypted/encrypted.tar.gz>. We also provide sample results of  
22  
23  
24  
25 400 merging 93 VCF files into either one VCF or one TPED file using our cluster-based schemas,  
26  
27  
28  
29 401 which are available on AWS S3 as  
30  
31  
32  
33 402 [https://s3.amazonaws.com/xsun316/sample\\_results/result.tar.gz](https://s3.amazonaws.com/xsun316/sample_results/result.tar.gz).

34  
35  
36  
37 403

#### 40 404 **Abbreviations**

41  
42  
43  
44 405 VCF: Variant Call Format; GWAS: Genome Wide Association Studies; WGS: Whole Genome  
45  
46  
47  
48 406 Sequencing; WES: whole exome sequencing; AWS: Amazon Web Service; CGDM: Collaborative  
49  
50  
51  
52 407 Genomic Data Model; SAM/BAM: Sequence/Binary Alignment/Map; RDD: Resilient Distributed  
53  
54  
55  
56 408 Dataset; DAG: Directed Acyclic Graph; EMR: Elastic-MapReduce; CAAPA: Consortium on Asthma  
57  
58  
59  
60 409 among African-ancestry Population in the Americas;  
61  
62  
63  
64  
65

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

410

411 **Consent for Publication**

412 Not applicable

413 **Competing Interests**

414 The authors declare they have no competing interests.

415 **Authors Contributions**

416 J.G. proposed the problem. X.S., F.W. initiated this project. X.S. designed and implemented the

417 CloudMerge project. X.S. drafted the manuscript. X.S., J.P., F.W. and Z.Q. revised the manuscript.

418 **Acknowledgements**

419 We are grateful for Duck, Alyssa Leann for helping on revising and rephrasing the manuscript.

420

421 **References**

- 422 1. Massie M, Nothaft F, Hartl C, Kozanitis C, Schumacher A, Joseph AD, et al. Adam: Genomics  
423 formats and processing patterns for cloud scale computing. University of California, Berkeley  
424 Technical Report, No UCB/EECS-2013. 2013;207.
- 425 2. Siretskiy A, Sundqvist T, Voznesenskiy M and Spjuth O. A quantitative assessment of the hadoop  
426 framework for analyzing massively parallel dna sequencing data. Gigascience. 2015;4 1:26.
- 427 3. Burren OS, Guo H and Wallace C. VSEAMS: a pipeline for variant set enrichment analysis using  
428 summary GWAS data identifies IKZF3, BATF and ESRRA as key transcription factors in type 1

- 1 429 diabetes. *Bioinformatics*. 2014;30 23:3342-8. doi:10.1093/bioinformatics/btu571.
- 2 430 4. Apache Hadoop. <http://hadoop.apache.org/>. Accessed 10 Oct 2017.
- 3
- 4 431 5. Dean J and Ghemawat S. Mapreduce: Simplified data processing on large clusters. *Commun Acm*.
- 5
- 6 432 2008;51 1:107-13. doi:Doi 10.1145/1327452.1327492.
- 7
- 8 433 6. Vora MN. Hadoop-HBase for large-scale data. In: *Computer science and network technology*
- 9
- 10 434 *(ICCSNT), 2011 international conference on 2011*, pp.601-5. IEEE.
- 11
- 12 435 7. Zaharia M, Chowdhury M, Das T, Dave A, Ma J, McCauley M, et al. Resilient distributed datasets:
- 13
- 14 436 A fault-tolerant abstraction for in-memory cluster computing. In: *Proceedings of the 9th*
- 15
- 16 437 *USENIX conference on Networked Systems Design and Implementation 2012*, pp.2-. USENIX
- 17
- 18 438 Association.
- 19
- 20 439 8. Schatz MC. CloudBurst: highly sensitive read mapping with MapReduce. *Bioinformatics*. 2009;25
- 21
- 22 440 11:1363-9. doi:10.1093/bioinformatics/btp236.
- 23
- 24 441 9. Langmead B, Schatz MC, Lin J, Pop M and Salzberg SL. Searching for SNPs with cloud computing.
- 25
- 26 442 *Genome Biol*. 2009;10 11:R134. doi:10.1186/gb-2009-10-11-r134.
- 27
- 28 443 10. Wang S, Mares MA and Guo YK. CGDM: collaborative genomic data model for molecular profiling
- 29
- 30 444 data using NoSQL. *Bioinformatics*. 2016;32 23:3654-60. doi:10.1093/bioinformatics/btw531.
- 31
- 32 445 11. AWS Genomics Guide.
- 33
- 34 446 [https://d0.awsstatic.com/Industries/HCLS/Resources/AWS\\_Genomics\\_WP.pdf](https://d0.awsstatic.com/Industries/HCLS/Resources/AWS_Genomics_WP.pdf). Accessed
- 35
- 36 447 10 Oct 2017.
- 37
- 38 448 12. Gruber K. Google for genomes. *Nature Research*, 2014.
- 39
- 40 449 13. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, et al. PLINK: a tool set for
- 41
- 42 450 whole-genome association and population-based linkage analyses. *Am J Hum Genet*. 2007;81
- 43
- 44 451 3:559-75. doi:10.1086/519795.
- 45
- 46 452 14. Mohammed EA, Far BH and Naugler C. Applications of the MapReduce programming framework
- 47
- 48 453 to clinical big data analysis: current landscape and future trends. *BioData Min*. 2014;7:22.
- 49
- 50 454 doi:10.1186/1756-0381-7-22.
- 51
- 52 455 15. Huang H, Tata S and Prill RJ. BlueSNP: R package for highly scalable genome-wide association
- 53
- 54 456 studies using Hadoop clusters. *Bioinformatics*. 2013;29 1:135-6.
- 55
- 56 457 doi:10.1093/bioinformatics/bts647.
- 57
- 58 458 16. White T. Hadoop: The definitive guide. " O'Reilly Media, Inc."; 2012.
- 59
- 60 459 17. Silberschatz A, Korth HF and Sudarshan S. *DatabaseSystem Concepts*. 2010.
- 61
- 62
- 63
- 64
- 65

- 1 460 18. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al. The variant call format  
2 461 and VCFtools. *Bioinformatics*. 2011;27 15:2156-8. doi:10.1093/bioinformatics/btr330.  
3  
4 462 19. Multiway-Merge Algorithm. [https://en.wikipedia.org/wiki/K-Way\\_Merge\\_Algorithms](https://en.wikipedia.org/wiki/K-Way_Merge_Algorithms).  
5  
6 463 Accessed 10 Oct 2017.  
7  
8 464 20. Chang F, Dean J, Ghemawat S, Hsieh WC, Wallach DA, Burrows M, et al. Bigtable: A distributed  
9  
10 465 storage system for structured data. *Acm T Comput Syst*. 2008;26 2 doi:Artn 4  
11  
12 466 10.1145/1365815.1365816.  
13  
14 467 21. Genomes Project C, Abecasis GR, Altshuler D, Auton A, Brooks LD, Durbin RM, et al. A map of  
15  
16 468 human genome variation from population-scale sequencing. *Nature*. 2010;467 7319:1061-73.  
17  
18 469 doi:10.1038/nature09534.  
19  
20 470 22. Mathias RA, Taub MA, Gignoux CR, Fu W, Musharoff S, O'Connor TD, et al. A continuum of  
21  
22 471 admixture in the Western Hemisphere revealed by the African Diaspora genome. *Nat*  
23  
24 472 *Commun*. 2016;7:12522. doi:10.1038/ncomms12522.  
25  
26 473 23. Kwon Y, Balazinska M, Howe B and Rolia J. A study of skew in mapreduce applications. *Open*  
27  
28 474 *Cirrus Summit*. 2011;11.  
29  
30 475 24. O'Neil P, Cheng E, Gawlick D and O'Neil E. The log-structured merge-tree (LSM-tree). *Acta*  
31  
32 476 *Informatica*. 1996;33 4:351-85.  
33  
34 477 25. Sedgewick R and Flajolet P. *An introduction to the analysis of algorithms*. Addison-Wesley; 2013.  
35  
36 478 26. Özsü MT and Valduriez P. *Principles of distributed database systems*. Springer Science & Business  
37  
38 479 *Media*; 2011.  
39  
40 480 27. Miner D and Shook A. *MapReduce Design Patterns: Building Effective Algorithms and Analytics*  
41  
42 481 *for Hadoop and Other Systems*. " O'Reilly Media, Inc."; 2012.  
43  
44 482 28. Chen L, Wang C, Qin ZS and Wu H. A novel statistical method for quantitative comparison of  
45  
46 483 multiple ChIP-seq datasets. *Bioinformatics*. 2015;31 12:1889-96.  
47  
48  
49 484

## 485 **Figure legends**

486 **Figure 1. Converting VCF files to TPED.** Left tables are input VCF files. Right table is the  
487 merged TPED file. Records are filtered out if their Filter value doesn't equal to 'PASS' (Pos

1 488 10147). Individual genotypes with the same genomic location that exist in any VCF file are  
2  
3  
4  
5 489 aggregated together on one row. The resulting TPED file thus has an inclusive set of sorted  
6  
7  
8  
9 490 genomic locations from all VCF files.  
10  
11  
12  
13 491  
14  
15  
16  
17 492 **Figure 2. The workflow chart of MapReduce schema.** It consists of two phases: In the first  
18  
19  
20  
21 493 phase, input VCF records are filtered, grouped around chromosomes into bins, and mapped into  
22  
23  
24  
25 494 key-value records. Two samplings are performed to generate partition lists of chromosomes. In the  
26  
27  
28  
29 495 second phase, parallel jobs of specified chromosomes are launched. Within each job, records from  
30  
31  
32  
33 496 corresponding bins are loaded, partitioned, sorted and merged by genomic locations before being  
34  
35  
36  
37 497 outputted as TPED files.  
38  
39  
40  
41 498  
42  
43  
44  
45 499 **Figure 3. The workflow chart of HBase schema.** The workflow is divided into three phases. The  
46  
47  
48  
49 500 first one is a sampling, filtering and mapping phase. A MapReduce job samples out VCF records  
50  
51  
52  
53 501 whose genomic positions are used as region boundaries when creating the HBase table. Only  
54  
55  
56  
57 502 qualified records are mapped as key-values and saved as Hadoop sequence files. The second phase  
58  
59  
60  
61 503 is HBase bulk loading in which a MapReduce job load and writes records outputted from the  
62  
63  
64  
65

1 504 previous phase, aggregating them into corresponding regional HFiles in the form of HBase's row  
2  
3  
4  
5 505 key and column families. Finished HFiles are moved into data folders on region servers. In the  
6  
7  
8  
9 506 third phase, we launch parallel scans over regions of the whole table to retrieve desired records  
10  
11  
12  
13 507 which are subsequently merged and exported as TPED files.  
14  
15  
16  
17 508

21 509 **Figure 4. The workflow chart of Spark schema.** It is a single Spark job consisting of three  
22  
23  
24  
25 510 stages. In the first stage, VCF records are loaded, filtered, and mapped to PairRDDs with keys of  
26  
27  
28  
29 511 genomic position and values of genotype. The sort-by-key shuffling spans across the first two  
30  
31  
32  
33 512 stages, sorting and grouping together records by keys. Then grouped records with the same key  
34  
35  
36  
37 513 are locally merged into one record in TPED format. Finally, merged records are exported as TPED  
38  
39  
40  
41 514 files.  
42  
43  
44  
45 515

48 516 **Figure 5. The scalability of clustered based schemas on input size.** Subfigures a, b and c refer  
49  
50  
51  
52 517 to MapReduce, HBase and Spark schemas respectively. As input file number increases from 10 to  
53  
54  
55  
56 518 93, the time cost of all three schemas with 12, 24 or 72 cores show at most a linear increasing  
57  
58  
59  
60 519 trend which suggests good scalabilities. The HBase schema with 12 cores has the largest increase  
61  
62  
63  
64  
65

1 520 from 375 to 2,751 seconds, or 7.3-fold.  
2  
3

4  
5 521  
6  
7

8  
9 522 **Figure 6. The scalability of cluster based schemas on available computing core number.**  
10

11  
12  
13 523 Subfigures a, b and c refer to MapReduce, HBase and Spark schemas respectively. In these  
14

15  
16  
17 524 experiments, the number of input files is fixed at 93. As core number increases from 12 to 72, the  
18

19  
20  
21 525 time cost of the three schemas decrease linearly until a plateau is reached where computing  
22

23  
24  
25 526 resources become excessive. The Spark schema shows lowest reduction of time cost from 1,699 to  
26

27  
28  
29 527 460 seconds, or 3.7-fold. These results suggest good scalabilities of these schemas on computing  
30

31  
32  
33 528 resources.  
34

35  
36  
37 529  
38

39  
40 530 **Figure 7. The performance anatomy of cluster-based schemas on increasing input size.** The  
41

42  
43  
44 531 number of cores in these experiments is fixed at 48. All phases of the three schemas show good  
45

46  
47  
48 532 scalabilities with input data size. **a)** MapReduce schema: The two MapReduce phases have a  
49

50  
51  
52 533 comparable time cost, increasing 3.0- and 2.2-fold respectively as input file number increases  
53

54  
55  
56 534 from 10 to 93. **b)** HBase schema: The time spent in each phase increases 2.0-, 2.7- and 2.2-fold  
57

58  
59  
60 535 respectively as input file number increases from 10 to 93. The bulk loading and exporting phases  
61  
62  
63  
64  
65

1 536 together take up more than 90% of total time expense. c) Spark schema: The time cost increases  
2  
3  
4  
5 537 3.1-, 3.0- and 3.4-fold respectively for the three stages as input file number increases from 10 to  
6  
7  
8  
9 538 93. Like the HBase schema, the first two stages of the Spark schema together account for more  
10  
11  
12  
13 539 than 90% of total time cost.  
14  
15  
16  
17 540  
18  
19  
20  
21 541 **Figure 8. Performance comparison among multiway-merge implementations and cluster-**  
22  
23  
24  
25 542 **based schemas:** Firstly, we compare of the performances of the three schemas with that of the  
26  
27  
28  
29 543 multiway-merge implementation. When input file number is 10, the time differences between  
30  
31  
32  
33 544 multiway-merge and our schemas are relatively small, ranging from 2- to 2.8-fold. As file number  
34  
35  
36  
37 545 increases to 93, the differences turn out to be more significant, ranging from 10.2- to 13.1-fold.  
38  
39  
40  
41 546 Secondly, we compare the performances among the three schemas which are comparable to each  
42  
43  
44  
45 547 other regardless of the input file number. MapReduce schema has best performance in merging 10  
46  
47  
48  
49 548 files; HBase schema performs best in merging 20, 40 and 60 files; Spark schema is fastest in  
50  
51  
52  
53 549 merging 93 files.  
54  
55  
56 550  
57  
58  
59  
60 551  
61  
62  
63  
64  
65

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

552 **Tables**

553 **Table 1. Performance comparisons of VCFTools versus MapReduce, HBase and Spark**

554 **schemas**

|                            | <b>VCFTools</b> | <b>MapReduce</b> | <b>HBase</b> | <b>Spark</b> |
|----------------------------|-----------------|------------------|--------------|--------------|
| <b>Time cost (seconds)</b> | 30,189          | 484              | 577          | 596          |
| <b>Fold (faster)</b>       | -               | 62.4             | 52.3         | 50.7         |

555 We compare the time cost of VCFTools and our three schemas using 72 cores on merging 40VCF  
556 files into one VCF file. VCFTools takes more than 30,000 seconds to finish. In contrast, all three  
557 schemas take less than 600 seconds to finish. MapReduce schema has the largest performance  
558 improvement which is about 62-fold.

559

560 **Table 2. Pros and Cons of MapReduce, HBase and Spark schemas**

| <b>Schemas</b>   | <b><i>Pros</i></b>  | <b><i>Cons</i></b>  |
|------------------|---|---|
| <b>MapReduce</b> | <ul style="list-style-type: none"><li>• Simple architecture and least overhead.</li><li>• Best parallelism for small input size (<math>\leq 20</math>).</li><li>• Good for one-time merging.</li><li>• Performance is stable.</li></ul> | <ul style="list-style-type: none"><li>• Merging is not incremental.</li></ul> |

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

|              |   |  |
|--------------|---|--|
| <b>HBase</b> | <ul style="list-style-type: none"> <li>• Good for intermediate input size (<math>\geq 20</math> and <math>\leq 100</math>).</li> <li>• Supports incremental merging.</li> <li>• Supports On-Line Analytical Processing (OLAP).</li> <li>• Best storage efficiency.</li> </ul> | <ul style="list-style-type: none"> <li>• Users must determine region number in advance.</li> <li>• Has most local I/O.</li> <li>• Complex performance tuning.</li> </ul>   |
| <b>Spark</b> | <ul style="list-style-type: none"> <li>• Good for large input size (<math>&gt; 100</math>).</li> <li>• Keeps intermediate results in memory and least local I/O.</li> <li>• Good for subsequent statistical analysis on merged results.</li> </ul>                            | <ul style="list-style-type: none"> <li>• Possibly weakened data locality during loading.</li> <li>• Slight unstable performance when computing resources exceeds needs of input size.</li> <li>• Actual execution plan is not transparent.</li> <li>• Complex performance tuning.</li> </ul> |

561 Each of the three distributed systems has its own specialties and limitations. As a result, the  
562 schemas running on them have different pros and cons, and application scenarios as listed above.

Figure1

[Click here to download Figure Figure1.pdf](#)VCF File  
1

| Chr | Pos   | Ref | Alt | Filter | ... | Genotype |
|-----|-------|-----|-----|--------|-----|----------|
| 1   | 10147 | A   | T   | q20    | ... | 1/0:43   |
| 1   | 10240 | T   | G   | PASS   | ... | 1/0:5    |
| ... | ...   | ... | ... | ...    | ... | ...      |
| Y   | 11590 | G   | C   | PASS   | ... | 0/0:10   |

VCF File  
2

| Chr | Pos   | Ref | Alt | Filter | ... | Genotype |
|-----|-------|-----|-----|--------|-----|----------|
| 1   | 10186 | G   | A   | PASS   | ... | 1/0:9    |
| 1   | 10240 | T   | G   | PASS   | ... | 1/1:11   |
| ... | ...   | ... | ... | ...    | ... | ...      |
| Y   | 11872 | G   | T   | PASS   | ... | 0/1:10   |

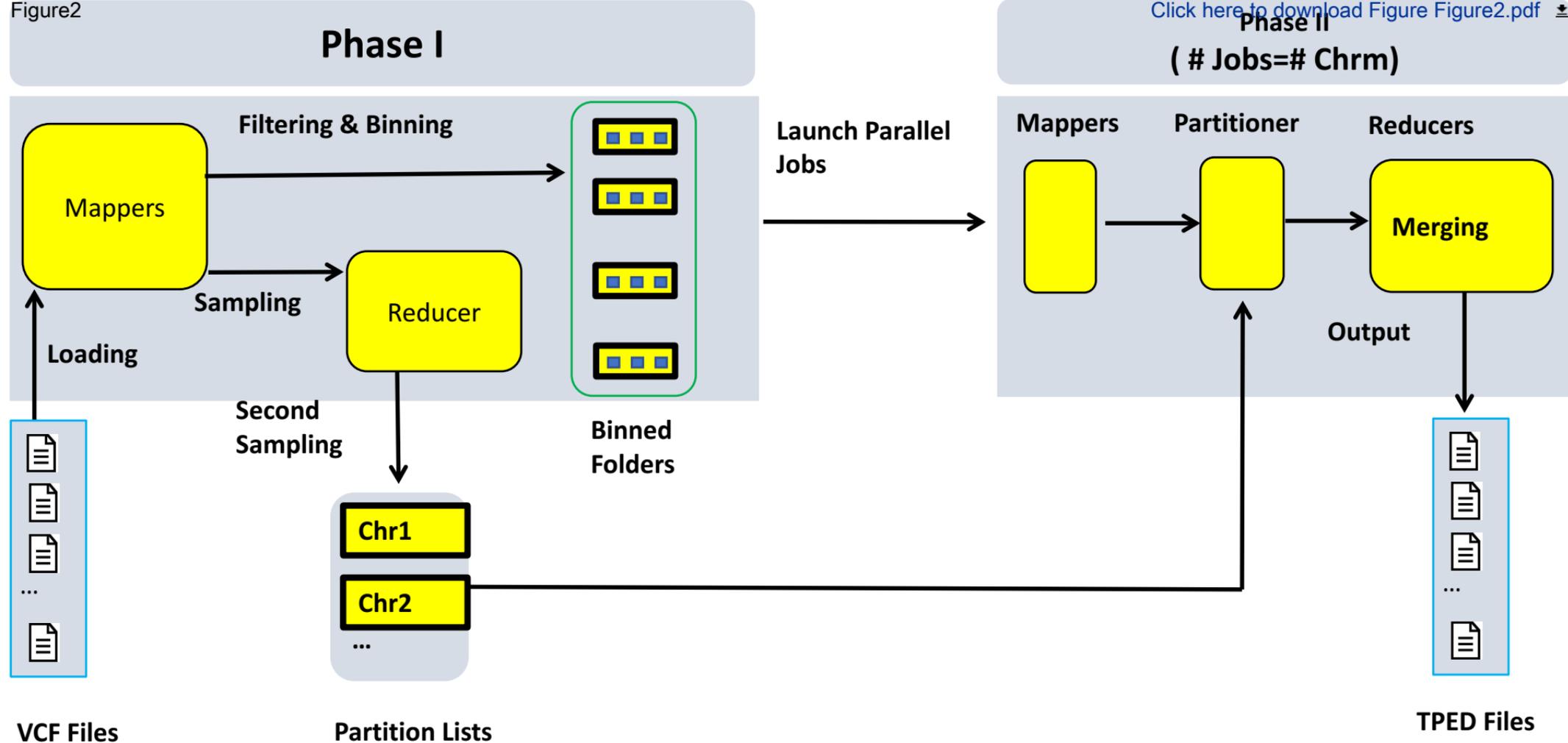


Genotypes

| Chr | Rs  | Distance | Pos   | Ind_1 | Ind_2 |
|-----|-----|----------|-------|-------|-------|
| 1   | .   | 0        | 10186 | G G   | G A   |
| 1   | .   | 0        | 10240 | T G   | G G   |
| ... | ... | ...      | ...   | ...   | ...   |
| Y   | .   | 0        | 11590 | G G   | G G   |
| Y   | .   | 0        | 11872 | G G   | G T   |

Merged TPED file

Figure2

[Click here to download Figure2.pdf](#)

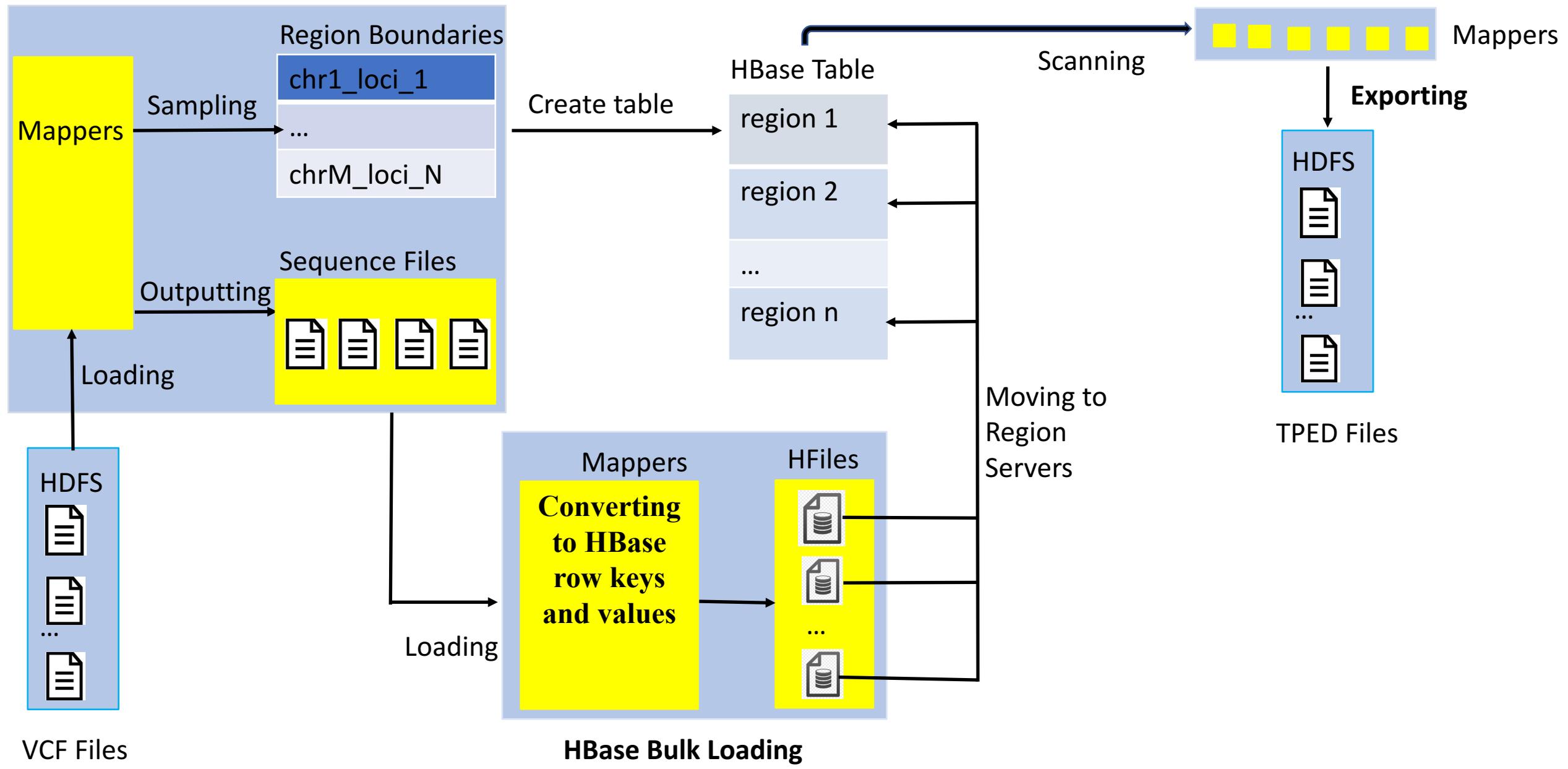
VCF Files

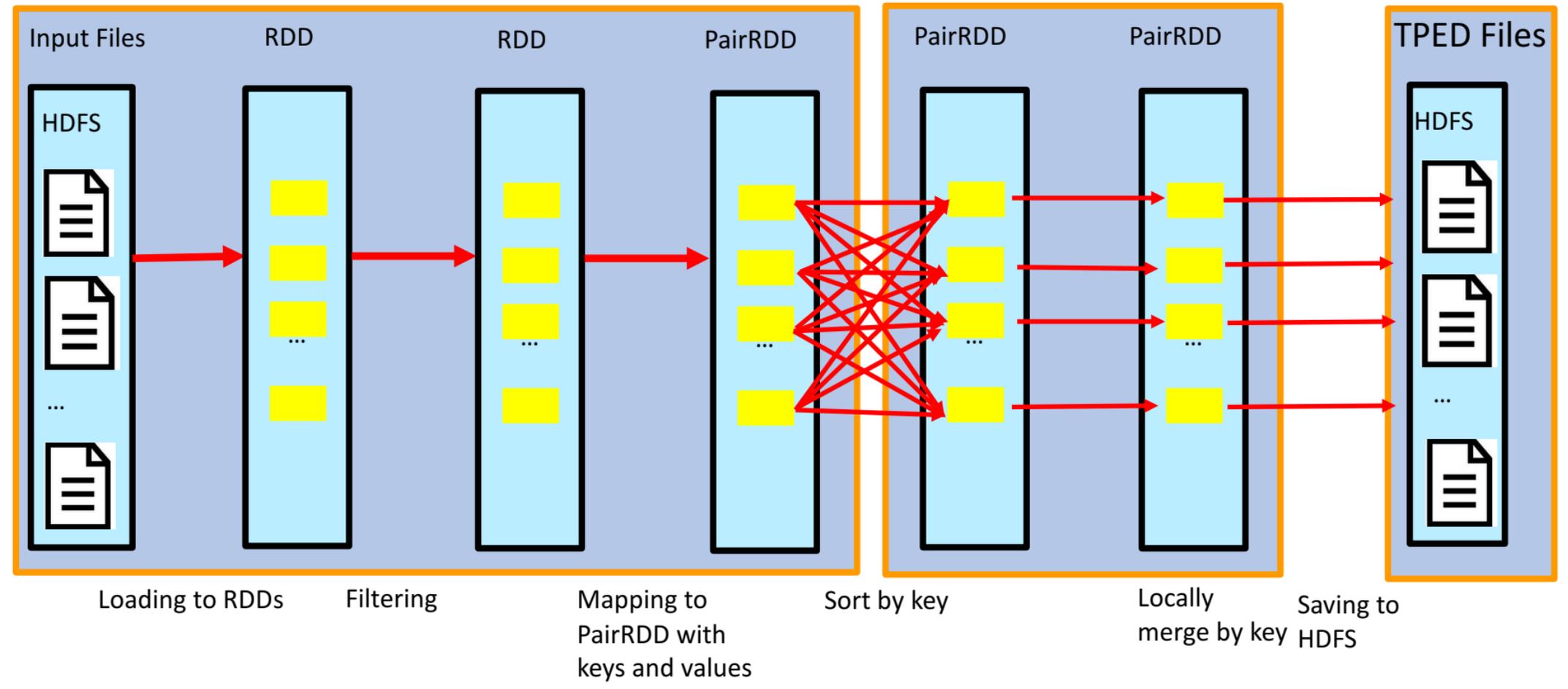
Partition Lists

TPED Files

Figure3

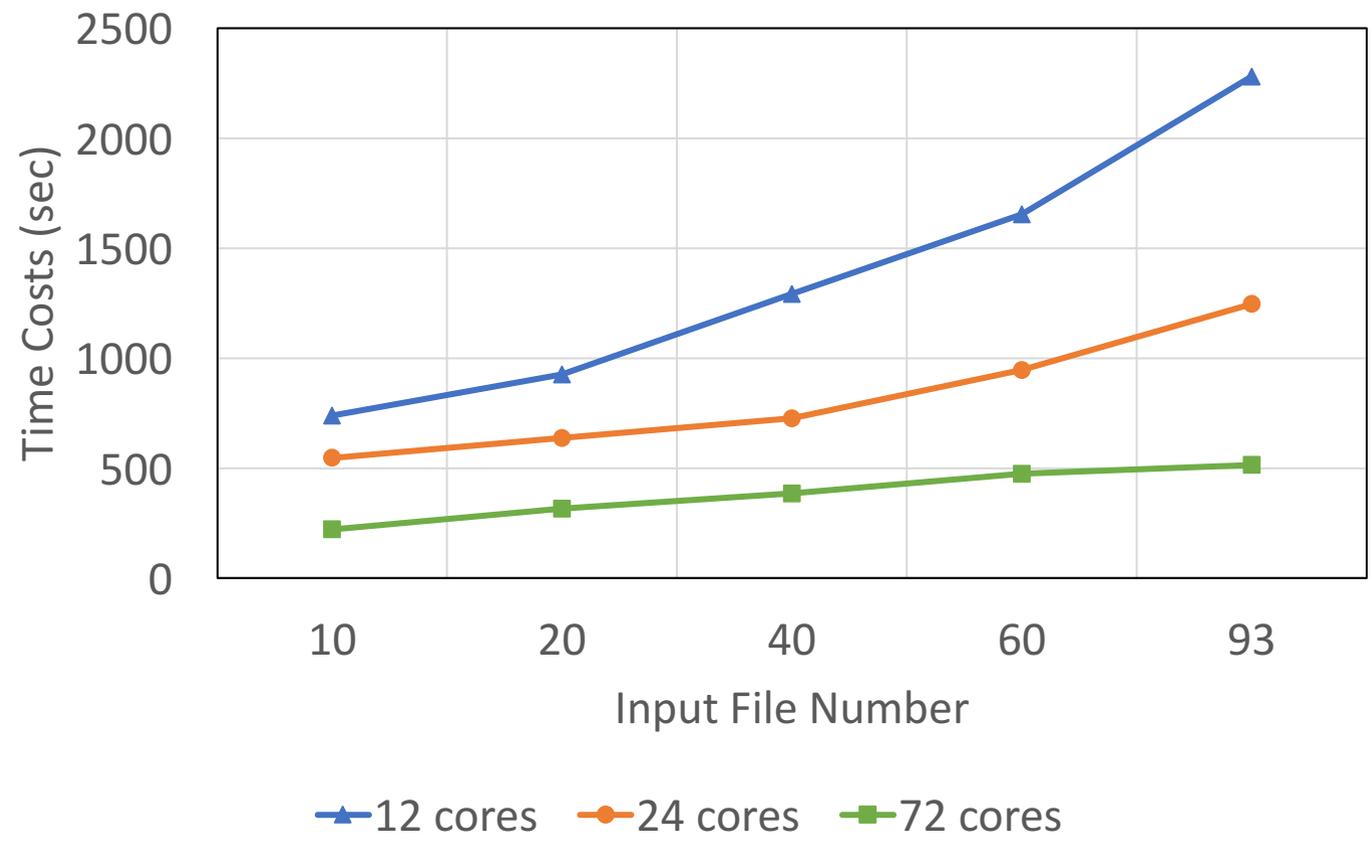
# Sampling, Mapping & Filtering





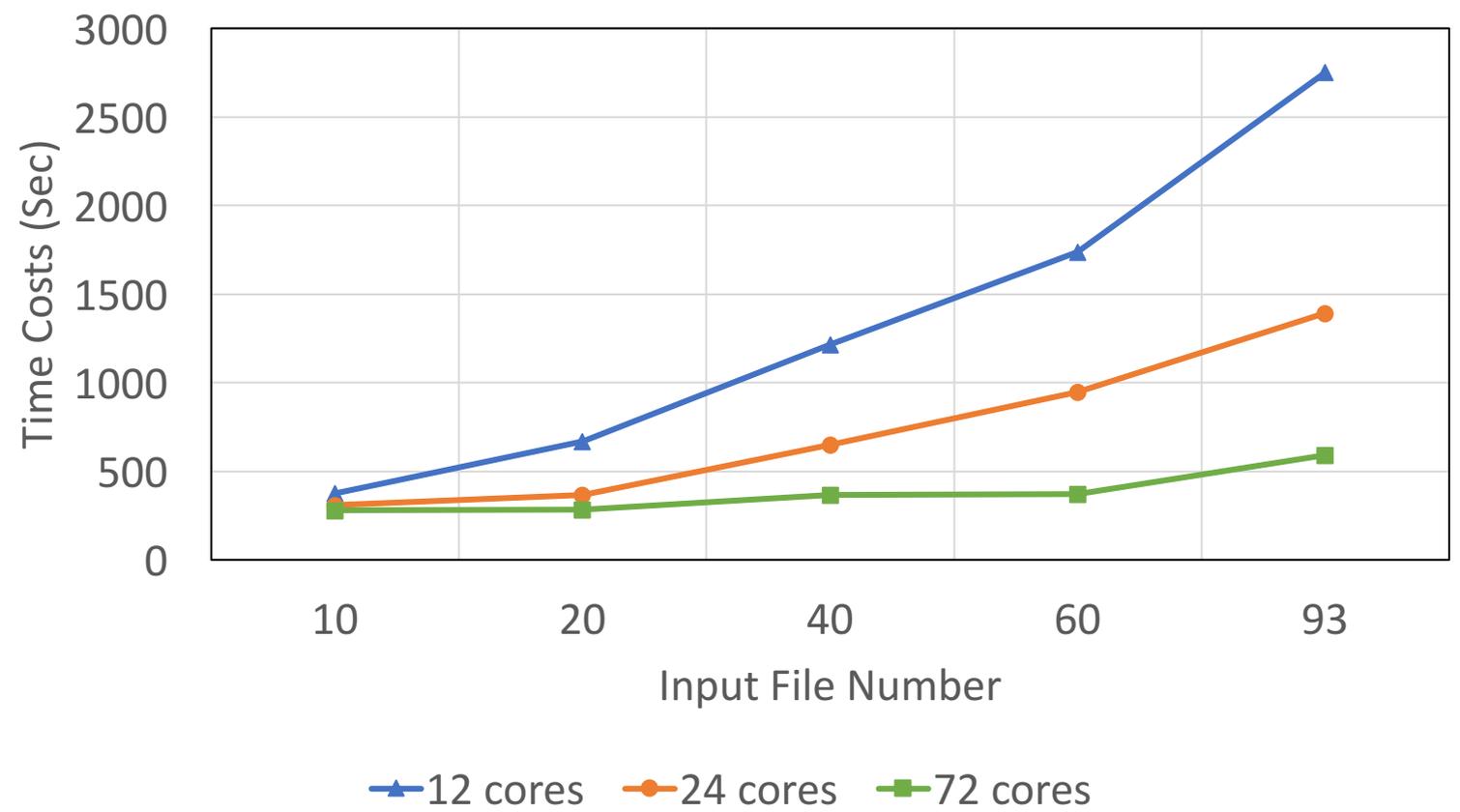
a)

MapReduce Schema



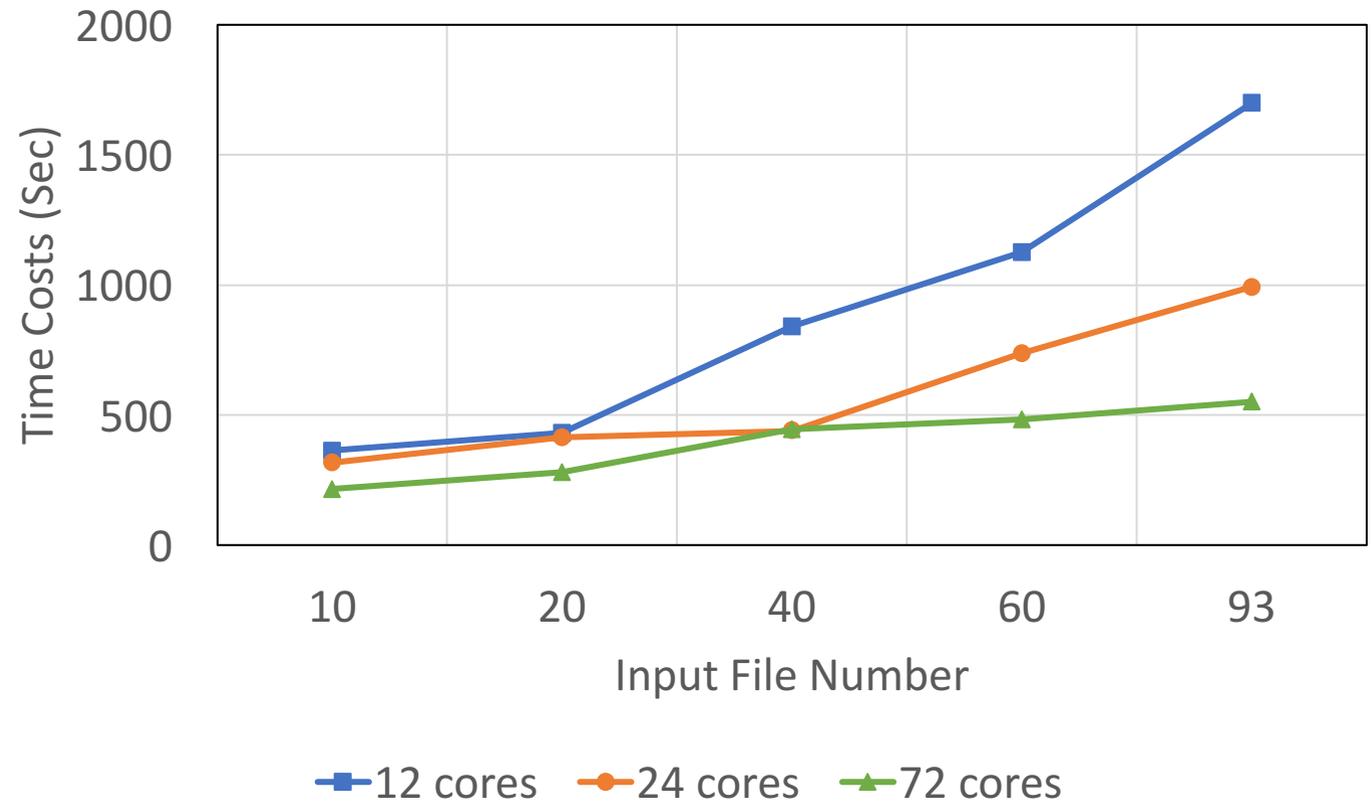
b)

HBase Schema

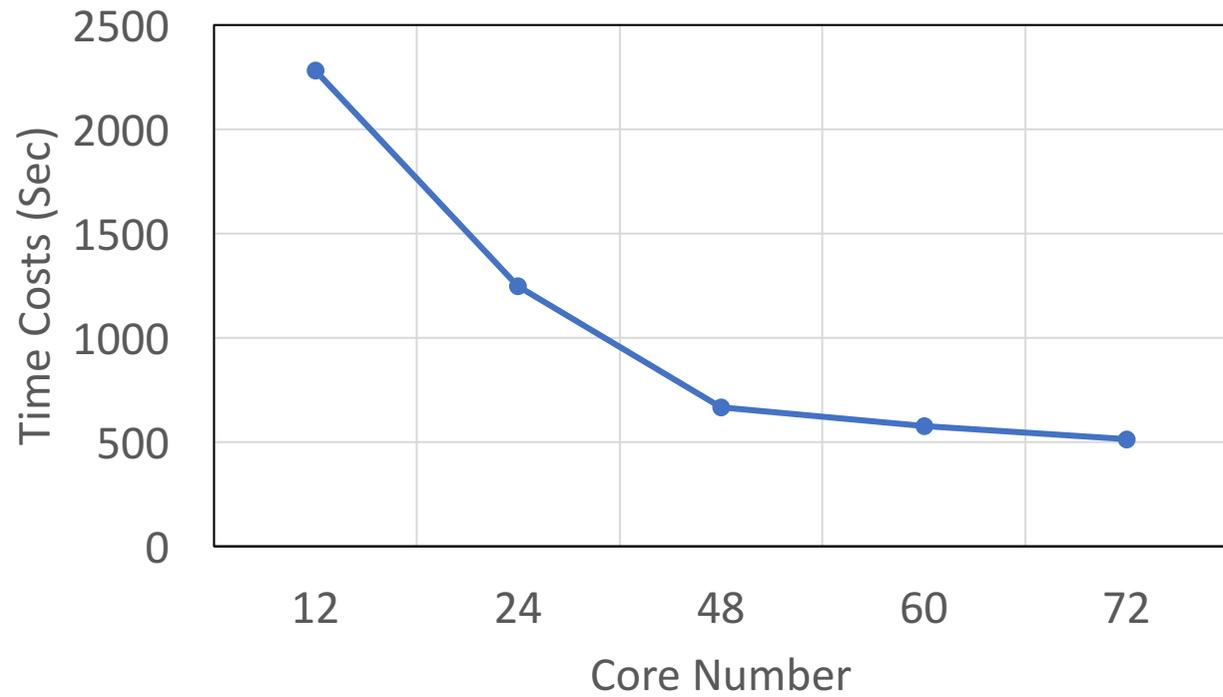


c)

Spark Schema

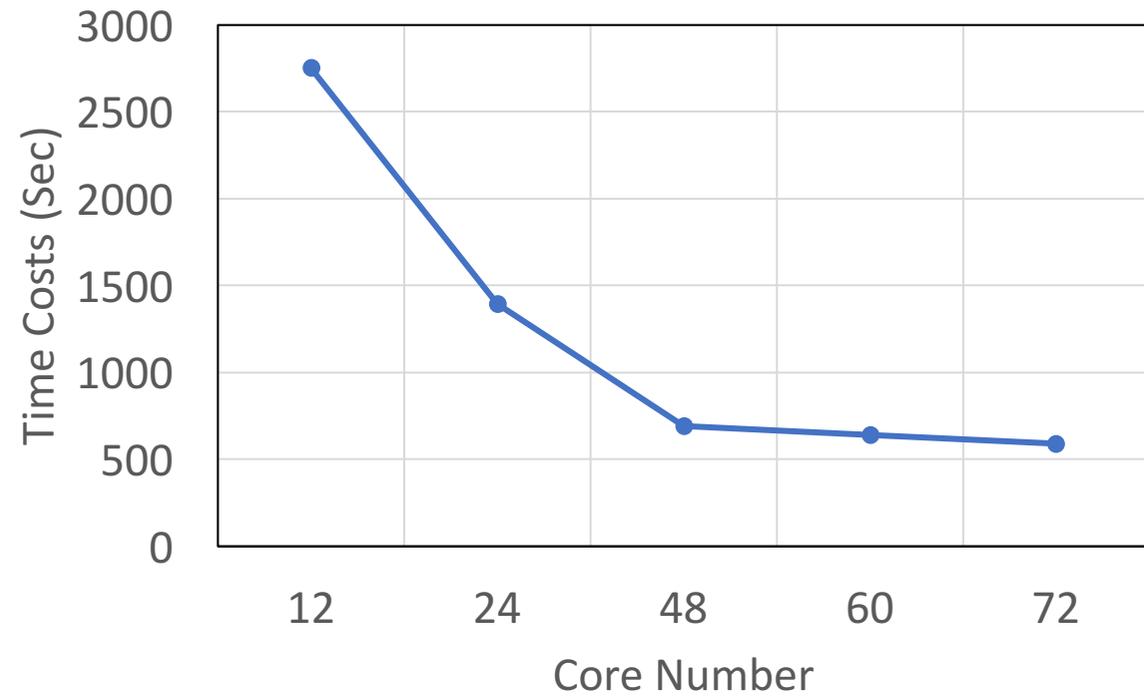


MapReduce Schema



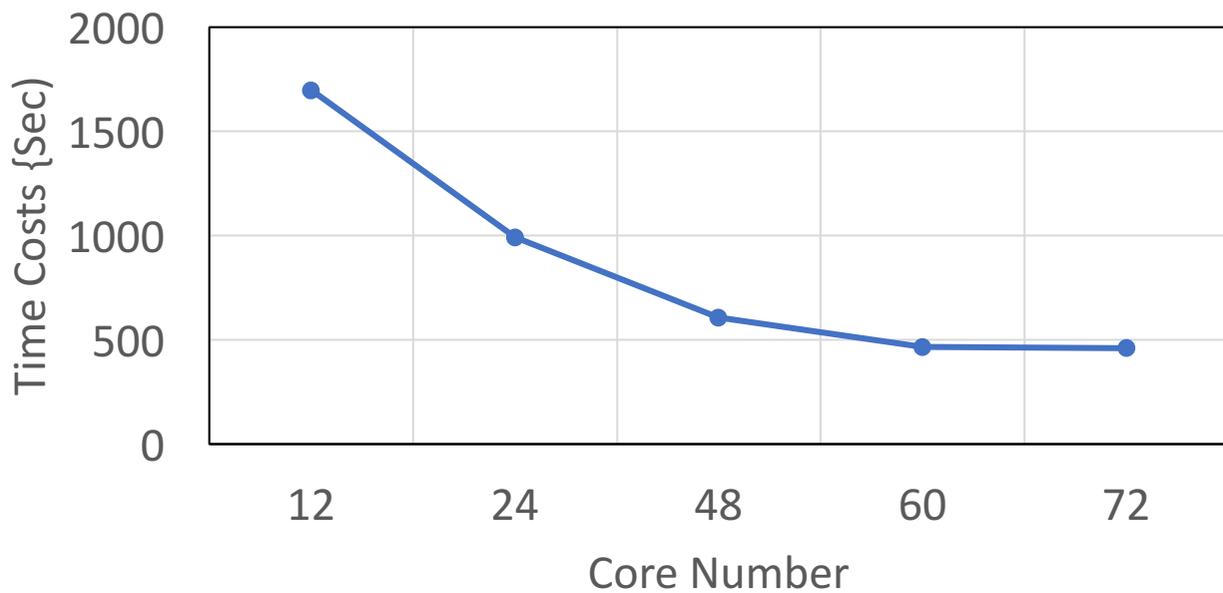
b)

HBase Schema



c)

Spark Schema



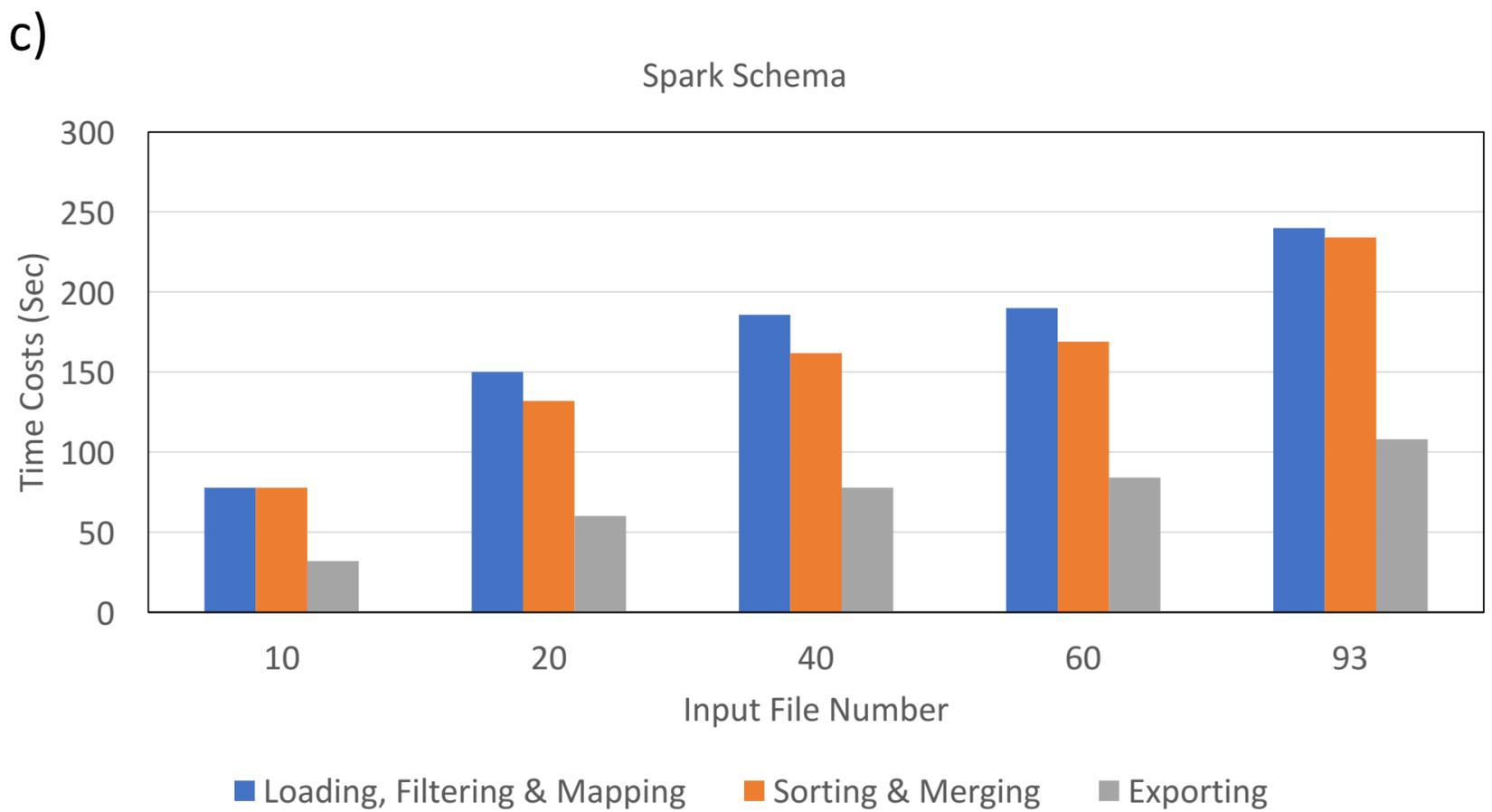
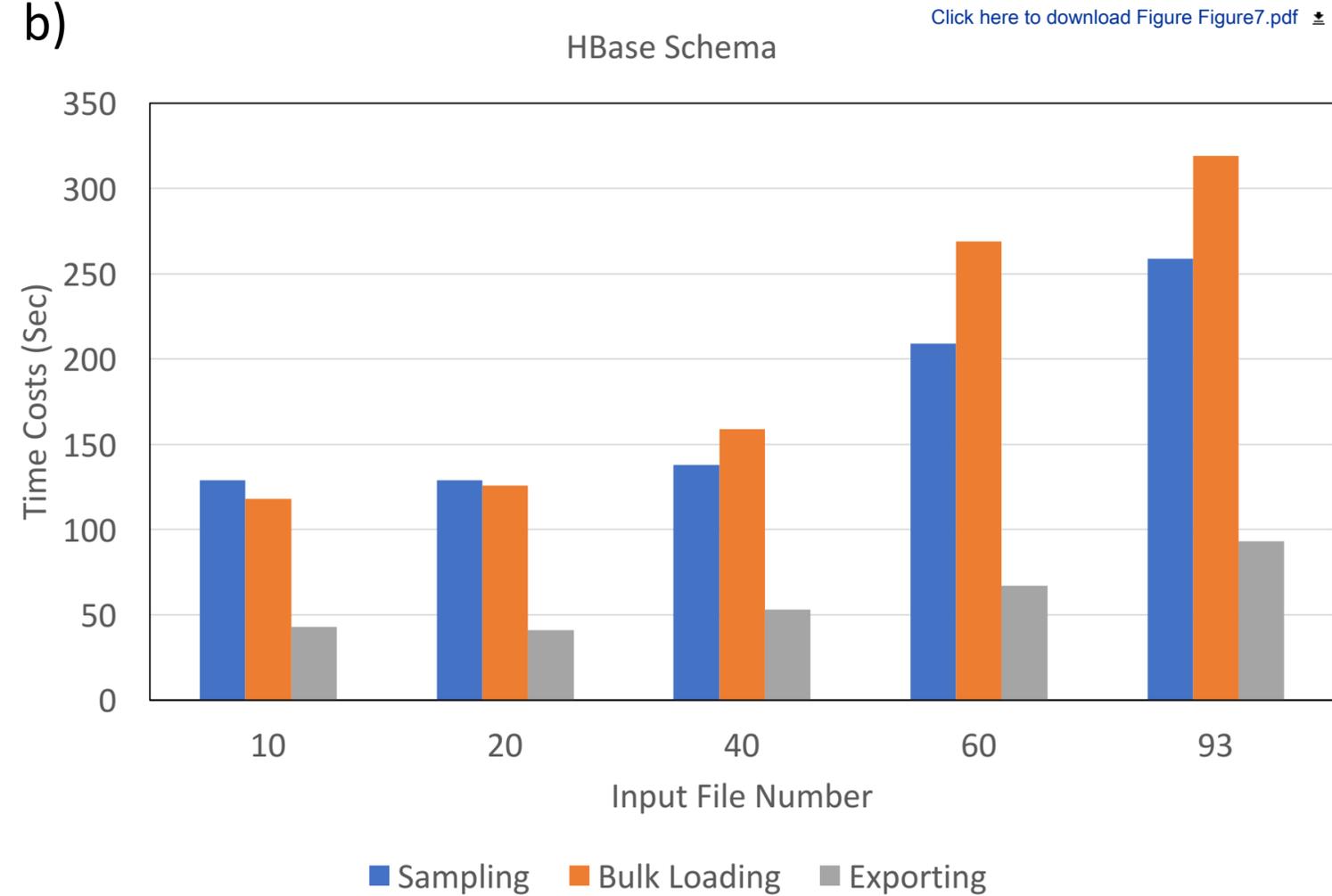
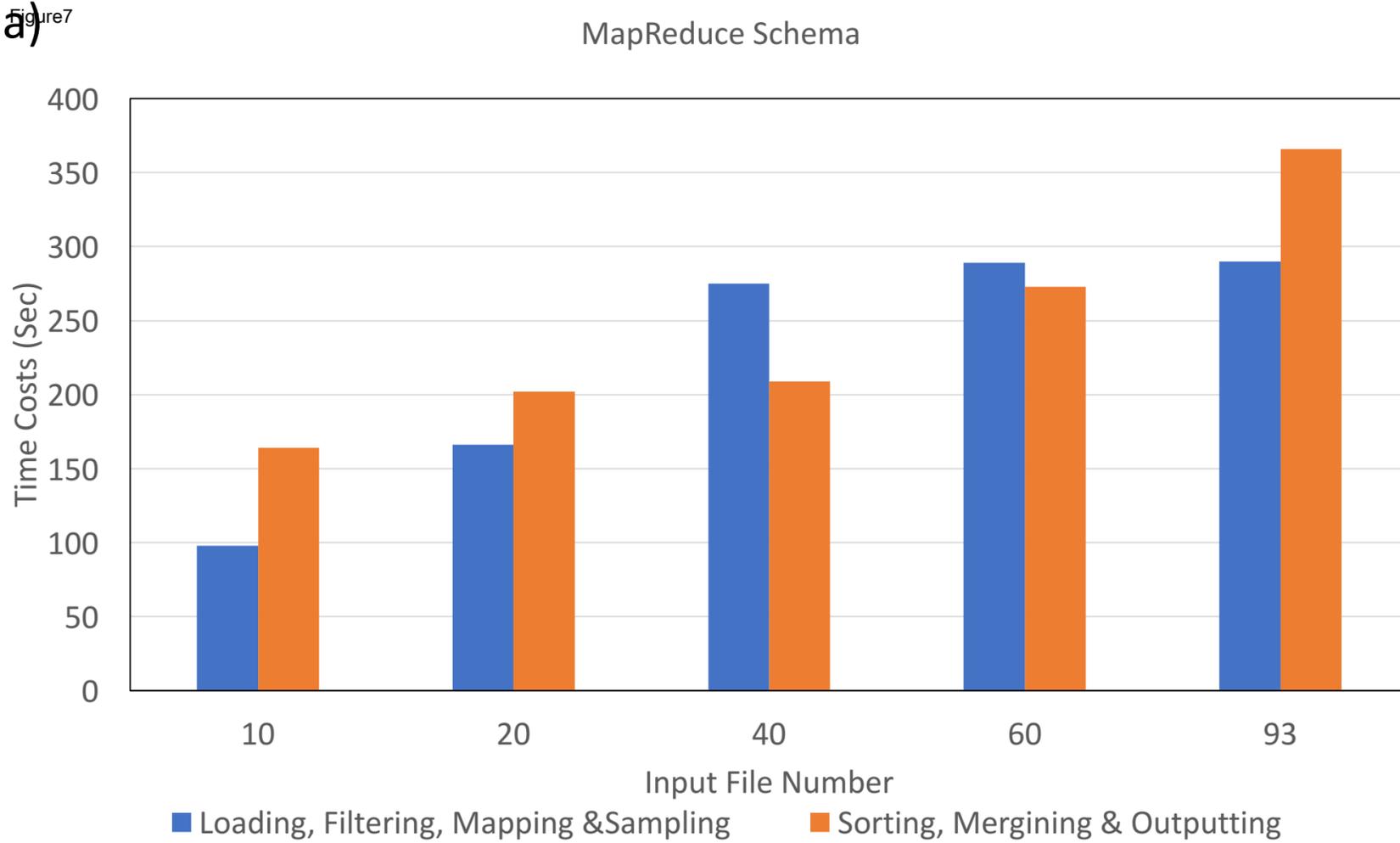
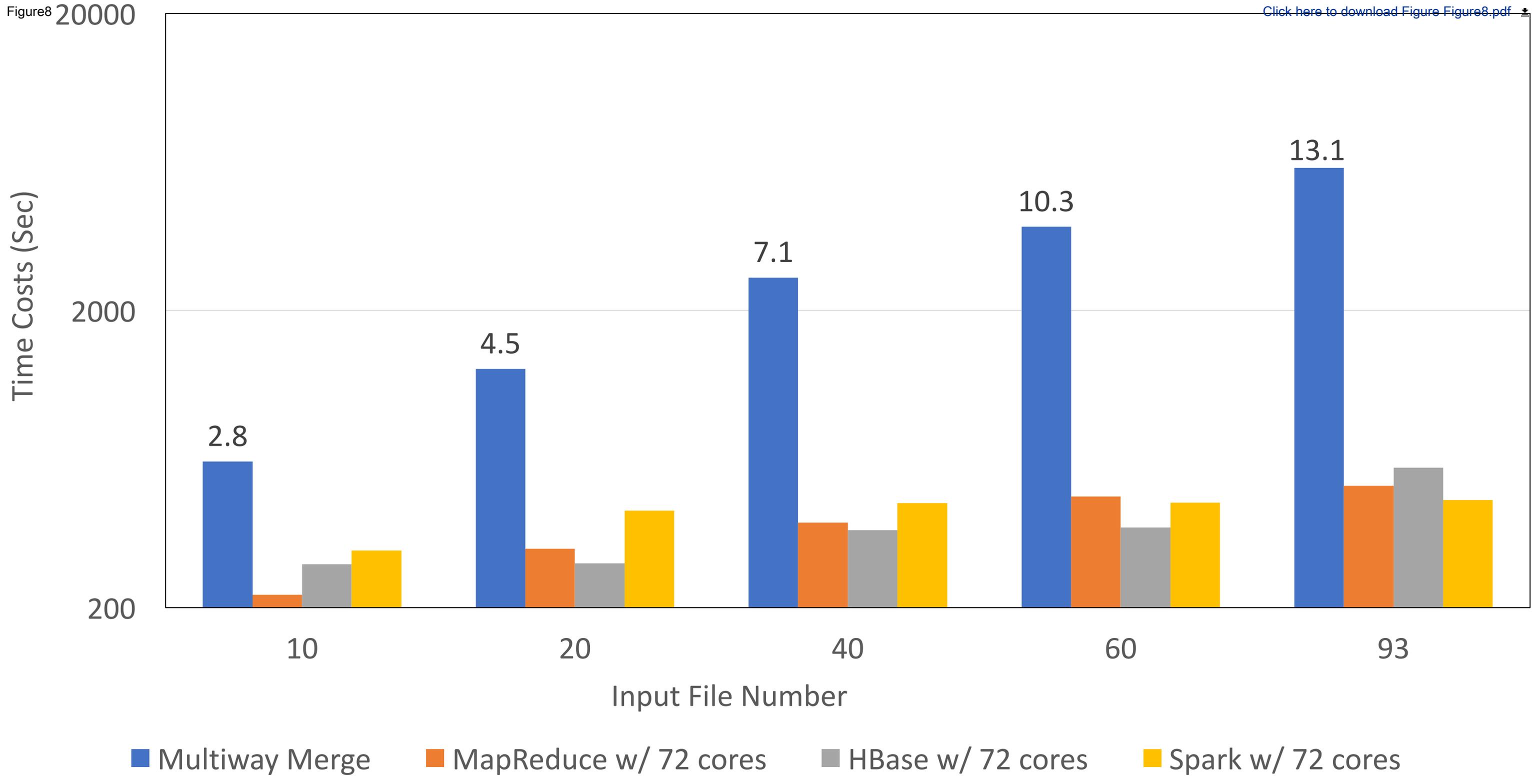
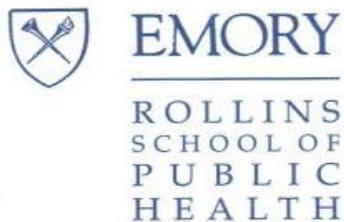


Figure8





Department of Biostatistics and Bioinformatics

October 12, 2017

Laurie Goodman, PhD  
Editor-in-Chief  
*GigaScience*

Dear Dr. Goodman:

On behalf of our colleague, Xiaobo Sun, Jingjing Gao and Peng Jin, I am pleased to submit a manuscript titled:

*Optimized Distributed Systems Achieve Significant Performance Improvement on Sorted Merging of Massive Omics Data*

for publication consideration in *GigaScience* as a Technical Note.

In the present manuscript, we describe three novel optimized schemas, running on Apache Hadoop, HBase and Spark respectively, for performing sorted merging on Omics data. Sorted merging is one of important data manipulation tools for Omics data, for example, the merged VCF or TPED files are required to perform statistical analysis in association studies on sequencing data. However, most existing tools for handling this task, such as VCFTools and PLINK, are running on single machine and implemented based on the multiway-merge algorithm. As a result, they suffer from the limitations of disk I/O and become very inefficient in face of large data size. The recent distributed systems offer an alternative solution. However, without optimized working schemas, naively using these systems does not lead to scalability.

In this study, we custom design optimized schemas for three Apache big data platforms. All three schemas are able to overcome the bottleneck problem by maintaining cluster's workload balance and achieving maximum parallelism. We have compared our schemas with VCFTools on merging 40 VCF files into a single one, as well as with a multiway-merge based implementation on merging up to 93 VCF files into a single TPED file. It turns out that even using a moderate sized cluster, we can archive speedup up to 62-fold compared to VCFTools. All three schemas show good scalability on both input size and number of cores, suggesting given enough computing resources we can guarantee the performance even in face of very large scale of data. Therefore our findings provide generalized scalable schemas for performing sorted merging on genetics and genomics data using these Apache distributed systems. Our schemas can be easily generalized to merge other types of Omics data such as CHIP-seq peak lists. **Therefore, we believe our schemas will have a high impact in the omics field as we enter the big data era.**

We hereby confirm that we do not have any potential competing interests and all authors have approved the manuscript for submission. We also confirm that the manuscript has not been submitted for publication elsewhere.

In light of the content, we suggested the following researchers as potential reviews for our manuscript: Dongxiao Zhu (Wayne State, email: [ct4442@wayne.edu](mailto:ct4442@wayne.edu)), Huanmei Wu (Indiana University–Purdue University Indianapolis, email: [hw9@iupui.edu](mailto:hw9@iupui.edu)), W. Jim Zheng (UT Health at Houston, email: [Wenjin.J.Zheng@uth.tmc.edu](mailto:Wenjin.J.Zheng@uth.tmc.edu)), Edmon Begoli (Oak Ridge National Laboratory, email: [begolie@ornl.gov](mailto:begolie@ornl.gov)), Ulf Leser (Humboldt-Universität zu Berlin, email: [leser@informatik.hu-berlin.de](mailto:leser@informatik.hu-berlin.de)).

Thank you very much for your kind editorial assistance.

Sincerely,



Zhaohui (Steve) Qin, Ph.D.  
Associate Professor  
Department of Biostatistics and Bioinformatics  
Emory University  
Atlanta, GA 30322

Fusheng Wang, Ph.D.  
Assistant Professor  
Department of Biomedical Informatics  
Department of Computer Science  
Stony Brook University  
2313D Computer Science,  
Stony Brook, NY 11794-8330