

GigaScience

Optimized Distributed Systems Achieve Significant Performance Improvement on Sorted Merging of Massive VCF Files

--Manuscript Draft--

Manuscript Number:	GIGA-D-17-00267R2	
Full Title:	Optimized Distributed Systems Achieve Significant Performance Improvement on Sorted Merging of Massive VCF Files	
Article Type:	Technical Note	
Funding Information:	Sandler Family Foundation	Dr Esteban Burchard
	American Asthma Foundation	Dr Esteban Burchard
	RWJF Amos Medical Faculty Development Program	Dr Esteban Burchard
	National Heart, Lung, and Blood Institute (R01HL104608)	Dr Kathleen Barnes
	National Institute of Neurological Disorders and Stroke (R01NS051630)	Dr Peng Jin
	National Institute of Neurological Disorders and Stroke (P01NS097206)	Dr Peng Jin
	National Institute of Neurological Disorders and Stroke (U54NS091859)	Dr Peng Jin
	National Science Foundation (ACI 1443054)	Dr Fusheng Wang
	National Science Foundation (IIS 1350885)	Dr Fusheng Wang
	National Heart, Lung, and Blood Institute (R01HL117004)	Dr Esteban Burchard
	National Heart, Lung, and Blood Institute (R01HL128439)	Dr Esteban Burchard
	National Heart, Lung, and Blood Institute (R01HL135156)	Dr Esteban Burchard
	National Heart, Lung, and Blood Institute (X01HL134589)	Dr Esteban Burchard
	National Institute of Environmental Health Sciences (R01ES015794)	Dr Esteban Burchard
	National Institute of Environmental Health Sciences (R21ES24844)	Dr Esteban Burchard
	National Institute on Minority Health and Health Disparities (P60MD006902)	Dr Esteban Burchard
	National Institute on Minority Health and Health Disparities (R01MD010443)	Dr Esteban Burchard
	National Institute on Minority Health and Health Disparities (RL5GM118984)	Dr Esteban Burchard
	Tobacco-Related Disease Research Program (24RT-0025)	Dr Esteban Burchard
Abstract:	<p>Background: Sorted merging of genomic data is a common data operation necessary in many sequencing-based studies. It involves sorting and merging genomic data from different subjects by their genomic locations. In particular, merging a large number of Variant Call Format (VCF) files is frequently required in large scale whole genome sequencing or whole exome sequencing projects. Traditional single machine based methods become increasingly inefficient when processing large numbers of VCF files due to the excessive computation time and I/O bottleneck. Distributed systems and</p>	

	<p>more recent cloud-based systems offer an attractive solution. However, carefully designed and optimized workflow patterns and execution plans (schemas) are required to take full advantage of the increased computing power while overcoming bottlenecks to achieve high performance.</p> <p>Findings: In this study, we custom design optimized schemas for three Apache big data platforms, Hadoop (MapReduce), HBase and Spark, to perform sorted merging of a large number of VCF files. 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 VCF files into either a single TPED or VCF file, which are benchmarked with the traditional single/parallel multiway-merge methods, message passing interface (MPI) based high performance computing (HPC) implementation and the popular VCFTools.</p> <p>Conclusions: Our experiments suggest all three schemas either deliver a significant improvement in efficiency or render much better strong and weak scalabilities over traditional methods. Our findings provide generalized scalable schemas for performing sorted merging on genetics and genomics data using these Apache distributed systems.</p>
Corresponding Author:	Zhaohui Qin UNITED STATES
Corresponding Author Secondary Information:	
Corresponding Author's Institution:	
Corresponding Author's Secondary Institution:	
First Author:	Xiaobo Sun
First Author Secondary Information:	
Order of Authors:	Xiaobo Sun
	Jingjing Gao
	Peng Jin
	Celeste Eng
	Esteban Burchard
	Terri Beaty
	Ingo Ruczinski
	Rasika Mathias
	Kathleen Barnes
	Fusheng Wang
	Zhaohui Qin
Order of Authors Secondary Information:	
Response to Reviewers:	<p>Major modifications in this revision are summarized below:</p> <ul style="list-style-type: none"> •Corrected grammar mistakes and minor wording issues. •Add members of CAAPA consortium to authorship. •Add "Ethics Approval and Consent to Participate" section. •Add "Funding" section. •Add CAAPA Consortium members to the "Authors Contributions" section. •Add CAAPA Consortium members in the "Acknowledgements" section •Add a complete list of CAAPA Consortium members and their affiliations at the end •Added our test datasets and codes to the GigaDB database.

- Added a citation of test datasets in GigaDB to the reference list, and cited it in appropriate places in the manuscript.
- Added all URLs as references in the bibliography, and cited them in the corresponding places in the manuscript.

Below we present itemized responses to all the comments, organized by editors and reviewers. The reviewers' comments are in bold. Our responses are in dark blue color.

Editor comments:

Before we hand over your manuscript to our production team:

- please go over the list of minor wording issues below, kindly provided by reviewer 1, and correct them in a revised submission.

Done

- please add a statement on ethical considerations to the manuscript. You are using encrypted VCFs, but as the original data set from reference [33] is under an authorized access scheme, if I understand correctly, I assume you needed IRB approval to access and work with the un-encrypted data for this project? Please also clarify in the manuscript whether this encrypted use of the subjects' data is covered by the consent they gave.

Done, please see "Ethics Approval and Consent to Participate" section in the manuscript.

Regarding your code and test data, one of our data curators will contact you shortly. Usually we host an archival copy of any code and test data in our repository GigaDB, which will be cited in the manuscript. Our data curators will discuss this with you.

Done

Please include a citation to any upcoming GigaDB dataset to your reference list (including the DOI link you will get from our data curators), and please cite this in the data availability section and elsewhere in the manuscript, where appropriate.

Please follow this example format for the reference:

[xx] Author1 N, Author2 N, AuthorX N. Supporting data for "Title of your manuscript". GigaScience Database. 2018. <http://dx.doi.org/xxxxxx>

(If you don't have a GigaDB doi at the time of resubmission, please leave the "dummy" version and we can exchange this for you.)

Please see reference 43 in the manuscript

Finally, a very minor point: Please include all URLs (except the "availability" section") as references in the bibliography, and cite them from the text rather than inserting them directly.

Done. On line 321:

The source codes are available at our GitHub website [35] (CloudMerge; RRID: SCR_016051).

Reviewer #1:

This paper has undergone substantial improvements since the original submission and the authors are to be commended on their efforts to address all the main issues raised in the initial review. I am satisfied that all of my concerns from the initial review have been adequately addressed, and I am happy to recommend that this paper is accepted for publication.

We are grateful to the reviewer for the comment.

I have a few minor comments below that the authors might wish to consider when editing the final version: "merging a large number of Variant Call Format (VCF) files are frequently encountered" -> "merging a large number of Variant Call Format (VCF) files is frequently encountered".

Done. See line 19.

"when processing hundreds or even thousands of VCF files" -> "when processing large volumes of VCF files".

Done. See line 22.

"The distributed systems and the more recent cloud-based systems" -> "Distributed systems and more recent cloud-based systems".

Done. See line 23.

"working flow" -> "workflow".

Done. See line 24.

"Apache Foundation has" -> "The Apache Foundation"

Done. See line 56.

"took advantage" -> "take advantage"

Thanks for brings this out. It is done, see line 66. In addition, we also make additional similar changes from past tense to current tense:

On line 64, "made" -> "make".

On line 68, "adopted" -> "adopts".

On line 69, "utilized" -> "utilizes".

Are two citations really needed for the "sorted full-outer-joining problem"? If it is well known, as the authors claim, then one citation should be sufficient.

Yes, we agree with the reviewer, and delete one reference: "28. Silberschatz A, Korth HF and Sudarshan S. DatabaseSystem Concepts. 2010."

"cumbersome" is probably the wrong word to describe the behaviour of PLINK and VCFTools on moderate numbers of input files. Cumbersome suggests that they are awkward or difficult to use, but really the problem is that their performance is unacceptable.

Yes, we agree with the reviewer.

On line 90:

"Currently, they are handled by software such as VCFTools [28] and PLINK, which become very cumbersome even in the face of a moderate number of VCF files."

Changed to:

Currently, they are handled by software such as VCFTools [28] and PLINK, which become considerably inefficient even in the face of a moderate number of VCF files.

"literally makes it sequential on writing" -> "makes it sequential on writing" (remove "literally", it is redundant) "

Done. See line 94.

and memory limitation" -> "and memory limits"

Done. See line 96.

"ideally fit" -> "is an ideal fit"

Done. See line 131.

"megabyte in size" -> "megabytes in size" (plural). Maybe use MB instead, to be consistent with the rest of the article using GB.

Done. See line 163.

"Key-value pairs" -> "key-value pairs" (capitalisation)

Done. See line 209.

It is not clear what this means "we only need to merge records from selected chromosomes of interest rather than from all of them". Can you please clarify?

We thank the reviewer 1 for bringing this up. Here we mean if we are only interested in the merged results of some specific chromosomes, say chr1-chr3, then we can just merge the records in the corresponding bins, instead of merging the records of all chromosomes. And on line 210, we have rephrased this sentence from "With this grouping, we only need to merge records from selected chromosomes of interest rather than from all of them." to "With this grouping, if SNPs of interest located in a few selected chromosomes only, we can choose to just merge records from these selected chromosomes rather than from all chromosomes."

delete: "which necessitates the adding of this phase"

Done. See line 261.

"finishing writings" -> "finishing writing" (not plural)

Done. See line 269.

It is likely that the HPC tests (namely the MPI version) would have performed better on a system with a high-performance file system such as GPFS or Lustre instead of NFS.

We totally agree with reviewer 1's opinion. Both Lustre and GPFS has better I/O scalability than NFS. And we expect our HPC benchmark would perform better using these file systems. However, I/O is not only the reason why the HPC benchmark does not scale well. Rather, the increasing in the number of merging rounds when increasing the number of input files is the main reason for decreasing efficiency. So we expect the scalability will improve to some extent but not too much when running it with the GPFS or Lustre file system. Another reason we choose NFS because we test our benchmark using StarCluster, which currently doesn't support either GPFS or Lustre.

Use GB for gigabytes instead of G.

Done.

Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
Experimental design and statistics	Yes
<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>	
<p>Resources</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 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>Yes</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>

1 **Optimized Distributed Systems Achieve Significant Performance**

2 **Improvement on Sorted Merging of Massive VCF Files**

3 **Xiaobo Sun¹, Jingjing Gao², Peng Jin³, Celeste Eng⁵, Esteban G. Burchard⁵, Terri H. Beaty⁶,**

4 **Ingo Ruczinski⁷, Rasika A. Mathias⁸, Kathleen C. Barnes⁴, Fusheng Wang^{9*}, Zhaohui Qin^{2,10*} ,**

5 **CAAPA consortium¹¹**

7 ¹Department of Computer Sciences, Emory University, Atlanta, GA 30322, USA.

8 ²Department of Medical Informatics, Emory University School of medicine, Atlanta, GA 30322, USA.

9 ³Department of Human Genetics, Emory University School of Medicine, Atlanta, GA 30322, USA.

10 ⁴Department of Medicine, University of Colorado Denver, Aurora, CO 80045

11 ⁵Department of Medicine, University of California, San Francisco, San Francisco, CA 94143 USA

12 ⁶Department of Epidemiology, Bloomberg School of Public Health, JHU, Baltimore, MD 21205 USA

13 ⁷Department of Biostatistics, Bloomberg School of Public Health, JHU, Baltimore, MD 21205 USA

14 ⁸Department of Medicine, Johns Hopkins University, Baltimore, MD 21224 USA

15 ⁹Department of Biomedical Informatics, Stony Brook University, Stony Brook, NY 11794, USA.

- 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 ¹⁰Department of Biostatistics, Emory University, Atlanta, GA 30322, USA.
- 17 ¹¹See end of the paper for a complete list of CAAPA members
- 18
- 19 X.S. Email: xsun28@emory.edu
- 20 J.G. Email: gaodoris@gmail.com
- 21 P.J. Email: peng_jin@emory.edu
- 22 C.E. Email: celeste.eng@ucsf.edu
- 23 E.G.B. Email: esteban.burchard@ucsf.edu
- 24 T.H.B. Email: tbeaty1@jhu.edu
- 25 I.R. Email: iruczin1@jhu.edu
- 26 R.A.M. Email: rmathias@jhmi.edu
- 27 K.C.B. Email: kathleen.barnes@ucdenver.edu
- 28 F.W. Email: fusheng.wang@stonybrook.edu
- 29 Z.Q. Email: zhaohui.qin@emory.edu
- 30 *Correspondence: zhaohui.qin@emory.edu, fusheng.wang@stonybrook.edu
- 31

1 32 **Abstract**

2
3
4
5 33 **Background:** Sorted merging of genomic data is a common data operation necessary in many
6
7
8
9 34 sequencing-based studies. It involves sorting and merging genomic data from different subjects by
10
11
12
13 35 their genomic locations. In particular, merging a large number of Variant Call Format (VCF) files
14
15
16
17 36 is frequently required in large scale whole genome sequencing or whole exome sequencing
18
19
20
21 37 projects. Traditional single machine based methods become increasingly inefficient when
22
23
24
25 38 processing large numbers of VCF files due to the excessive computation time and I/O bottleneck.
26
27
28
29 39 Distributed systems and more recent cloud-based systems offer an attractive solution. However,
30
31
32
33 40 carefully designed and optimized workflow patterns and execution plans (schemas) are required to
34
35
36
37 41 take full advantage of the increased computing power while overcoming bottlenecks to achieve
38
39
40
41 42 high performance.

42
43
44
45 44 **Findings:** In this study, we custom design optimized schemas for three Apache big data platforms,
46
47
48
49 45 Hadoop (MapReduce), HBase and Spark, to perform sorted merging of a large number of VCF
50
51
52
53 46 files. These schemas all adopt the divide-and-conquer strategy to split the merging job into
54
55
56
57 47 sequential phases/stages consisting of subtasks which are conquered in an ordered, parallel and
58
59
60
61
62
63
64
65

1 48 bottleneck-free way. In two illustrating examples, we test the performance of our schemas on
2
3
4
5 49 merging multiple VCF files into either a single TPED or VCF file, which are benchmarked with
6
7
8
9 50 the traditional single/parallel multiway-merge methods, message passing interface (MPI) based
10
11
12
13 51 high performance computing (HPC) implementation and the popular VCFTools.
14
15
16
17 52

18
19
20
21 53 **Conclusions:** Our experiments suggest all three schemas either deliver a significant improvement
22
23
24
25 54 in efficiency or render much better strong and weak scalabilities over traditional methods. Our
26
27
28
29 55 findings provide generalized scalable schemas for performing sorted merging on genetics and
30
31
32
33 56 genomics data using these Apache distributed systems.
34
35

36 57 **Keywords:** Sorted merging, whole genome sequencing, MapReduce, Hadoop, HBase, Spark.
37
38
39
40
41 58

42 43 44 59 **Findings**

45 46 47 48 60 **Introduction**

49
50
51
52 61 With the rapid development of high-throughput biotechnologies, genetic studies have entered the
53
54
55
56 62 Big Data era. Studies like Genome Wide Association Studies (GWASs), Whole Genome
57
58
59
60 63 Sequencing (WGS) and whole exome sequencing (WES) studies have produced massive amounts
61
62
63
64
65

1 64 of data. The ability to efficiently manage and process such massive data becomes increasingly
2
3
4
5 65 important for successful large scale genetics studies [1-3]. Single machine based methods are
6
7
8
9 66 inefficient when processing such big data due to the prohibitive computation time, I/O bottleneck,
10
11
12
13 67 as well as CPU and memory limitations. Traditional HPC techniques based on MPI/OpenMP also
14
15
16
17 68 suffer from limitations such as not allowing addition of computing nodes at runtime, shortage of a
18
19
20
21 69 fault-tolerant and high available file system, inflexibility of customizing the computing
22
23
24
25 70 environment without administrator permission of a cluster [3, 4]. It becomes increasingly
26
27
28
29 71 attractive for investigators to take advantage of more powerful distributed computing resources or
30
31
32
33 72 the cloud to perform data processing and analyses [3, 5]. The Apache Foundation has been a
34
35
36
37 73 leading force in this endeavor, and has developed multiple platforms and systems including
38
39
40
41 74 Hadoop [6, 7], HBase [8] and Spark [9]. All these three Apache platforms have gained substantial
42
43
44
45 75 popularity in recent years, and have been endorsed and supported by major vendors such as
46
47
48
49 76 Amazon Web Services (AWS).

50
51
52 77
53
54
55
56

57 78 In bioinformatics, researchers have already started to embrace Apache distributed systems to
58
59
60
61
62
63
64
65

1 79 manage and process large amounts of high throughput ‘-omics’ data. For example, the Cancer
2
3
4
5 80 Genome Atlas project makes use of the Hadoop framework to split genome data into chunks
6
7
8
9 81 distributed over the cluster for parallel processing[3, 10]. The CloudBurst [11], Seal [12], Hadoop-
10
11
12
13 82 BAM [13] and Crossbow software [14] take advantage of the Hadoop framework to accelerate
14
15
16
17 83 sequencing read mapping, aligning and manipulations as well as SNP calling. The Collaborative
18
19
20
21 84 Genomic Data Model (CGDM) [15] adopts HBase to boost the querying speed for the main
22
23
24
25 85 classes of queries on genomic databases. MetaSpark [16] utilizes Spark’s distributed data set to
26
27
28
29 86 recruit large scale of metagenomics reads to reference genomes, achieves better scalability and
30
31
32
33 87 sensitivity than single-machine based programs [17]. Industry cloud computing vendors such as
34
35
36
37 88 Amazon [18] and Google [19] are also beginning to provide specialized environments to ease
38
39
40
41 89 genomics data processing in the cloud.

42
43
44
45 90

46
47
48
49 91 Although numerous Apache cluster-based applications have already been developed for
50
51
52
53 92 processing and analyzing large scale genomics data including ADAM [1], VariantSpark [20],
54
55
56
57 93 SparkSeq [21], Halvade [22], SeqHBase [23] among others, we believe there are still many
58
59
60
61
62
63
64
65

1 94 opportunities in biomedical data analyses to take advantage of distributed systems as the scale and
2
3
4
5 95 scope of data become larger and more complex. A particular example is sorted merging, which is a
6
7
8
9 96 ubiquitous operation in processing genetics and genomics data. As an example, in WGS, variants
10
11
12
13 97 identified from individuals are often called and stored in separate Variant Call Format (VCF) files.
14
15
16
17 98 Eventually these VCF files need to be merged (into a VCF or TPED file) as required by
18
19
20
21 99 downstream analysis tools such as PLINK [24] and BlueSNP [25, 26]. Either a VCF or TPED file
22
23
24
25 100 requires the data to be sorted by their genomic locations, thus these tasks are equivalent to the
26
27
28
29 101 well-known sorted full-outer-joining problem [27]. Currently, they are handled by software such
30
31
32
33 102 as VCFTools [28] and PLINK, which become considerably inefficient even in the face of a
34
35
36
37 103 moderate number of VCF files. The main reason is that these tools adopt the multiway-merge-like
38
39
40
41 104 method [29] with a priority queue as the underlying data structure to ensure the correct output
42
43
44
45 105 order. Although such a method only requires one round of read through of the input files, a key
46
47
48
49 106 deficiency is that it can only have one consumer access items from the data queue, which makes it
50
51
52
53 107 sequential upon writing. This problem cannot be eliminated even if the multiway-merging is
54
55
56
57 108 implemented as parallel processes due to I/O saturation, workload imbalance among computing
58
59
60
61 109 units, and memory limits. Therefore, these single-machine based tools are inefficient and time-

1 110 consuming when handling large datasets.
2
3
4
5
6 111
7
8
9
10 112 In this study, we use the case of sorted-merging multiple VCF files to demonstrate the benefits of
11
12
13
14 113 using Apache distributed platforms. However, simply running sorted merging on such distributed
15
16
17
18 114 systems runs into problems of bottlenecks, hotspots and unordered results commonly seen in
19
20
21
22 115 parallel computations. Rather, we believe working schemas custom designed for each specific
23
24
25
26 116 distributed platform are required to unleash their full potential. To overcome the limitations of
27
28
29
30 117 single-machine, traditional parallel/distributed, and simple Apache distributed system based
31
32
33
34 118 methods, we propose and implement three schemas running on Hadoop, Spark and HBase
35
36
37
38 119 respectively. We choose these three platforms because they represent cloud distributed systems
39
40
41
42 120 providing data partitioning based parallelism with distributed storage, data partitioning based
43
44
45
46 121 parallelism with in-memory based processing, and high dimensional tables like distributed
47
48
49
50 122 storage, respectively. Hadoop [6] is the open source implementation of MapReduce [7] based on
51
52
53
54 123 parallel key-value processing technique, and has the advantage of transparency and simplicity.
55
56
57
58 124 HBase [8] is a data warehousing platform which adopts Google's BigTable data storing structure
59
60
61
62
63
64
65

1 125 [30] to achieve high efficiency in storing and reading/writing large scale of sparse data. Spark [9]
2
3
4
5 126 introduces the concept of Resilient Distributed Dataset (RDD) and Directed Acyclic Graph (DAG)
6
7
8
9 127 execution to parallel key-value processing, thus enabling fast, robust and repetitive in-memory
10
11
12
13 128 data manipulations. Specifically, our schemas involve dividing the job into multiple phases
14
15
16
17 129 corresponding to tasks of loading, mapping, filtering, sampling, partitioning, shuffling, merging
18
19
20
21 130 and outputting. Within each phase, data and tasks are evenly distributed across the cluster,
22
23
24
25 131 enabling processing large scale of data in a parallel and scalable manner, which in turn improves
26
27
28
29 132 both speed and scalability.

30
31
32
33 133

34 35 36 37 134 **Methods**

38 39 40 41 135 **Overview**

42
43
44
45 136 The benefits of using these three Apache distributed platforms to perform sorted merging are four-
46
47
48
49 137 fold when compared to using the multiway-merge method [29], a relational database based
50
51
52
53 138 approach, or a HPC framework. First, with genomic locations as keys and genotypes as values, it
54
55
56
57 139 is readily transformed into the key-value model in which all three platforms offer a rich set of
58
59
60
61
62
63
64
65

1 140 parallel operations. Second, data in VCF files are semi-structured. This type of data is an ideal fit
2
3
4
5 141 for the three platforms which allow defining the schema during data loading, avoiding the
6
7
8
9 142 preprocessing of raw data into a rigid schema as in a relational database. Third, all these
10
11
12
13 143 platforms provide built-in efficient task coordination, high fault tolerance, data availability and
14
15
16
17 144 locality which are absent in the traditional HPC framework. Fourth, the merged results are directly
18
19
20
21 145 saved onto a distributed file system such as HDFS or Amazon S3 which can be directly used for
22
23
24
25 146 subsequent cluster-based GWAS or WGS analytical tools such as BlueSNP.
26
27
28
29 147
30
31
32
33 148 Despite these advantages, simply performing sorted merging on these Apache distributed systems
34
35
36
37 149 will not deliver the expected results for the following reasons. First, it can lead to globally
38
39
40
41 150 unsorted results. Hash-based shuffling of input data is the default mechanism for distributing data
42
43
44
45 151 to parallel working units in the system. However, shuffling will lead to globally unsorted results.
46
47
48
49 152 Second, bottlenecks and hotspots can happen during the processing in the cluster. Bypassing the
50
51
52
53 153 hashing based shuffling can lead to unbalanced workloads across the cluster, result in straggling
54
55
56
57 154 computing units which become bottlenecks for response time. In addition, for parallel loading of
58
59
60 155 presorted data into HBase, data being loaded from all the loading tasks access the same node
61
62
63
64
65

1 156 simultaneously while other nodes may be idling, creating an I/O hotspot. Third, sampling costs
2
3
4
5 157 could become prohibitive. Although Hadoop provides a built-in utility named *total-order-merging*
6
7
8
9 158 [27] to achieve both workload balance and global order, it involves transferring to and sampling
10
11
12
13 159 all the data on a single node. The communication costs over the network and disk I/O can be
14
15
16
17 160 prohibitive when data size becomes very large. In the following sections, we will illustrate how
18
19
20
21 161 our custom designed schemas are able to overcome these limitations in detail.

22
23
24
25 162

26 27 28 163 **Data Formats and Operations**

29
30
31
32 164 In a typical WGS experiment, data analysis often starts from individual genotype files in the VCF
33
34
35
36 165 format [31]. A VCF file contains data arranged into a table consisting of eight mandatory fields
37
38
39
40 166 including chromosome (CHROM), the genomic coordinate of the start of the variant (POS), the
41
42
43
44 167 reference allele (REF), a comma separated list of alternate alleles (ALT), among others. In our
45
46
47
48 168 experiments, we use a dataset consisting of the VCF files of 186 individuals [32] generated from
49
50
51
52 169 Illumina's BaseSpace software (Left tables in Figure 1). Each VCF file has around 4-5 million
53
54
55
56 170 rows, each row contains information on one of the individual's genomic variants. Each VCF file is
57
58
59
60 171 about 300 MB in size. In an attempt to protect privacy of study subjects, we apply the following
61
62
63
64
65

1 172 strategy to conceal their real genetic variant information contained in the VCF files: we first
2
3
4
5 173 transform each original genomic location by multiplying it with an undisclosed constant real
6
7
8
9 174 number, taking the floor integer of the result, and then add another undisclosed constant integer
10
11
12
13 175 number.
14
15
16
17 176
18
19
20
21 177 It is common that multiple VCF files need to be merged into a single TPED file for analysis tools
22
23
24
25 178 such as PLINK. A TPED file resembles a big table, aggregating genotypes of all individuals under
26
27
28
29 179 investigation by genomic locations (right table in Figure 1). The merging follows several rules.
30
31
32
33 180 First, each record is associated with a data quality value in the FILTER column, which records the
34
35
36
37 181 status of this genomic position passing all filters. Usually only qualified records with a “PASS”
38
39
40
41 182 filter value are retained. Second, genotypes in VCF files are stored in the form of allele values,
42
43
44
45 183 where 0 stands for the reference allele, 1 stands for the first mutant allele, 2 stands for the second
46
47
48
49 184 mutant allele, and so on. Allele values must be translated into corresponding types of nucleotides
50
51
52
53 185 in the TPED file. Third, all individuals need to have a genotype for genomic locations appearing
54
55
56
57 186 in at least one VCF file. The default genotype for a missing value is a pair of homozygous
58
59
60
61 187 reference alleles. The merging of multiple VCF files into a single VCF file follows the rules as:
62
63
64
65

1 188 First, the ALT and INFO columns of a genomic location in the merged file are set as the
2
3
4
5 189 concatenated values of the corresponding columns on that location from all input files with
6
7
8
9 190 duplicated values removed. Second, the QUAL column of a genomic location in the merged file is
10
11
12
13 191 set as a weight-averaged quality value of all individuals on that location. Third, a genomic
14
15
16
17 192 location is kept only when it appears in at least one input file and has a FILTER column value of
18
19
20
21 193 "PASS". Fourth, if an individual does not have allele values on a genomic location in the input
22
23
24
25 194 file, their missing allele values are designated as "." in the merged file.
26
27
28
29 195
30
31
32 196 For our Apache cluster-based schemas, the merging of multiple VCF files into a single TPED file
33
34
35
36
37 197 and the merging of multiple VCF files into a single VCF file differ only in the value contents of
38
39
40
41 198 the key-value pairs, so they should have the same scalability property. Although we implement
42
43
44
45 199 the applications of both merging types using our Apache cluster-based schemas, which are
46
47
48
49 200 available on our project website, we focused our experiments on the merging of multiple VCF
50
51
52
53 201 files into a single TPED file and only evaluate the execution speed of the merging of multiple
54
55
56
57 202 VCF files into a single VCF file with VCFTools as the benchmark.
58
59
60 203
61
62
63
64
65

1 204 **MapReduce (Hadoop) Schema**

2
3
4
5 205 This schema is built on Hadoop's underlying model MapReduce and running on Hadoop clusters.

6
7
8
9 206 MapReduce [7] is a parallel computing model based on a *split-apply-combine* strategy for data

10
11
12
13 207 analysis, in which data are mapped to key-values for splitting (mapping), shuffling and combining

14
15
16
17 208 (reducing) for final results. We use Apache Hadoop-2.7 as the system for our implementation. Our

18
19
20
21 209 optimized schema consists of two MapReduce phases, as shown in Figure 2 (the pseudocodes are

22
23
24
25 210 shown in Figure S1).

26
27
28
29 211

30
31
32 212 **1) First MapReduce phase.**

33
34
35
36 213 Raw data are loaded from HDFS into parallel mappers to perform the following tasks: First,

37
38
39
40 214 unqualified data are filtered out and qualified ones are mapped to key-value pairs. The mapper

41
42
43
44 215 output key is the genomic location and output value is the genotype and individual ID. Second,

45
46
47
48 216 key-value pairs are grouped together by chromosomes and temporarily saved as compressed

49
50
51
52 217 Hadoop sequence files [33] for faster I/O in the second MapReduce phase. With this grouping, if

53
54
55
56 218 SNPs of interest located in a few selected chromosomes only, we can choose to just merge records

57
58
59
60 219 from these selected chromosomes rather than from all chromosomes. Meanwhile, these records are

1 220 sampled to explore their distribution profile of keys along chromosomes to determine boundaries.
2
3
4
5 221 The boundaries are determined so there is an approximately equal number of records within each
6
7
8
9 222 segment. Because all records falling in the same segment will be assigned to the same reducer in a
10
11
12
13 223 later phase, boundaries calculated in this way ensure the workload of each reducer is balanced.
14
15
16
17 224 There are two rounds of samplings. The first one happens in each mapper with a pre-specified
18
19
20
21 225 sampling rate, which in our case is set to be 0.0001. Sampled records are then separated and
22
23
24
25 226 distributed to different reducers in this phase by chromosomes, where they are sampled again with
26
27
28
29 227 a rate equal to the reciprocal of the number of input files. This second sampling effectively limits
30
31
32
33 228 the number of final sampled records even in the face of a very large number of input files. Because
34
35
36
37 229 the number of reducers instantiated in the second phase equals the number of boundaries, which in
38
39
40
41 230 turn is decided by the number of sampled records, we can therefore avoid launching unnecessary
42
43
44
45 231 reducers thus minimizing task overheads.

46
47
48 232

51
52 233 **2) *Second MapReduce phase.***
53
54
55

56 234 In this phase, multiple parallel MapReduce jobs are created, one for each chromosome, to handle
57
58
59
60 235 all the records in sequence files generated from the first phase. Within each job, a partitioner
61
62
63
64
65

1 236 redirects records to the appropriate reducer by referring to the splitting boundaries from the
2
3
4
5 237 previous phase, so records falling in between the same pair of boundaries are aggregated together.
6
7
8
9 238 Finally, each reducer sorts and merges aggregated records by genomic locations before saving
10
11
12
13 239 them to a TPED file. In this way, globally sorted merging can be fulfilled.
14
15
16
17 240

21 241 **HBase Schema**

22
23
24 242 HBase [8] is a column-oriented database where data are grouped into column families and split
25
26
27
28 243 horizontally into regions spreading across the cluster. With this data storing structure, it supports
29
30
31
32 244 efficient sequential reading and writing of large-scale data as well as fast random data accessing.
33
34
35
36 245 Also, HBase is storage efficient because it can remember null values without saving them on disk.
37
38
39
40 246 These features make HBase an ideal platform for managing large, sparse data with relatively low
41
42
43
44 247 latency which naturally fits the sorted merging case. We use the HBase-1.3 as the system for our
45
46
47
48 248 implementation. As shown in Figure 3, our optimized HBase schema is divided into three phases
49
50
51
52 249 as discussed next (refer to Figure S2 for pseudocodes).
53
54
55
56 250

60 251 *1) Sampling phase*

1 252 The main challenge of HBase is to avoid computational hotspots in the cluster which can happen
2
3
4
5 253 when it starts loading a table from a single region hosted by a single node. Therefore, we need to
6
7
8
9 254 presplit the table into regions of approximately equal size before loading. The sampling phase is
10
11
12
13 255 introduced to determine reasonable presplitting regional boundaries. The total number of regions
14
15
16
17 256 is set to half of the number of input files so the size of each region is approximately 1GB.
18
19
20
21 257 Meanwhile, mappers of this phase also save qualified records as compressed Hadoop sequence
22
23
24
25 258 files on HDFS which are used as inputs in the next phase. In addition, filtering and key-value
26
27
28
29 259 mapping also take place in this phase.

30
31
32
33 260

34 35 36 261 **2) Bulk loading phase**

37
38
39
40 262 Even when the table has been presplit evenly, the hotspot problem of loading sorted inputs can
41
42
43
44 263 still emerge because sorted records are loaded sequentially, and at any instant they still access the
45
46
47
48 264 same region and server. During the bulk loading, the key and value of each record produced from
49
50
51
52 265 the previous phase is converted into HBase's binary row-key and column-value respectively, and
53
54
55
56 266 saved into a HFile, HBase's native storage format. The row-key here is in the form of
57
58
59
60 267 chromosome-genomic location, and column-value refers to reference allele, individual ID and
61
62
63
64
65

1 268 genotype. The bulk loading populates each HFile with records falling in the same pair of presplit
2
3
4
5 269 regional boundaries. Because HFiles are written simultaneously by parallel mappers/reducers, all
6
7
8
9 270 working nodes are actively involved and the regional hotspot is thus circumvented. Upon finishing
10
11
12
13 271 writing, the HBase can readily load HFiles in parallel into the table by simply moving them into
14
15
16
17 272 local HBase storage folders. This procedure is therefore at least an order of magnitude faster than
18
19
20
21 273 the normal loading in which data are loaded sequentially via HBase servers' I/O routines. The
22
23
24
25 274 order of records in the table is guaranteed because they are internally sorted by writing reducers
26
27
28
29 275 and HBase's Log-Structured Merge-tree [34]. It worth mentioning that VCF records are always
30
31
32
33 276 sparse, thus HBase is very storage-efficient.

34
35
36
37 277

38
39
40 278 **3) *Exporting phase***

41
42
43
44 279 A scan of a specified genomic window is performed on the table. It involves launching parallel
45
46
47
48 280 mappers each receiving records from a single HBase region, filling in missing genotypes,
49
50
51
52 281 concatenating records with the same row-key, and outputting final results into TPED files.

53
54
55
56 282

57
58
59
60 283 **Spark Schema**

1 284 Spark [9] is a distributed engine built upon the ideas of MapReduce and RDD. It can save
2
3
4
5 285 intermediate results in the form of RDD in memory, and perform computations on them. Also, its
6
7
8
9 286 computations are lazily evaluated, which means the execution plan can be optimized to include as
10
11
12
13 287 many computational steps as possible. As a result, it is ideal for iterative computations such as
14
15
16
17 288 sorted merging. We implement our optimized Spark schema on Spark-2.1. It has three stages
18
19
20
21 289 which we describe below and present in Figure 4 (refer to Figure S3 for pseudocodes).
22
23
24

25 290

27
28
29 291 **1) *RDD preprocessing stage***

30
31
32 292 This stage involves loading raw data as RDDs, filtering, and mapping RDDs to paired-RDDs with
33
34
35
36 293 keys (chromosome and genomic position) and values (reference allele, sample ID and genotype).
37
38
39

40 294 This stage ends with a sorting-by-key action which extends to the next stage.
41
42
43

44 295

46
47
48 296 **2) *Sorting and merging stage***

49
50
51
52 297 The sort-by-key shuffling repartitions and sorts PairRDD records so records with the same key
53
54
55
56 298 are aggregated together, which are then merged into the TPED format and converted back to RDD
57
58
59
60 299 records for outputting. However, Spark's native family of group-by-key functions for merging
61
62
63
64
65

1 300 should not be used here because their default partitioner is hash-based and different from the
2
3
4
5 301 range-based partitioner used by previous sort-by-key function. Consequently, the merged results
6
7
8
9 302 would be reshuffled into an unsorted status. We therefore optimize the merging to bypass these
10
11
12
13 303 functions so merging can be performed locally without data reshuffling to ensure both order and
14
15
16
17 304 high speed.

20
21 305

22
23
24
25 306 **3) *Exporting stage***

26
27
28
29 307 In this stage, merged RDD records are saved as TPED files on HDFS.

30
31
32
33 308

34
35
36
37 309 Execution parallelism has an important impact on the performance. To maximize performance, the
38
39
40
41 310 number of parallel tasks is set to be the number of input files. In this way, data locality is
42
43
44
45 311 maximized and each task is assigned a proper amount of work. In addition, unlike using
46
47
48
49 312 MapReduce or HBase, when performing sorting by keys, no explicit sampling is needed because
50
51
52
53 313 Spark keeps track of the number of records before determining repartition boundaries.

54
55
56 314

57
58
59
60 315 **Parallel Multiway-Merge and MPI-based High Performance Computing Implementations**

1 316 For most bioinformatics researchers, their daily working environment is still traditional in-house
2
3
4
5 317 HPC clusters or stand-alone powerful servers (with cores ≥ 16 and memory $\geq 200\text{GB}$) rather than
6
7
8
9 318 heterogeneous cloud-based clusters. Therefore, we also implement a parallel multiway-merge
10
11
12
13 319 program running on a single machine and a MPI-based (mpi4py v3.0) “single program, multiple
14
15
16
17 320 data (SPMD)” program running on a HPC cluster as benchmarks. The source codes are available
18
19
20
21 321 at our GitHub website [35] (CloudMerge; RRID: SCR_016051). We choose to implement
22
23
24
25 322 multiway-merge, because many existing bioinformatics tools, including VCFTools and PLINK,
26
27
28
29 323 adopt it as the underlying algorithm for sorted merging. Multiway-merge is highly efficient on
30
31
32
33 324 single machine as it requires only one scan of sorted input files, so it can theoretically run at the
34
35
36
37 325 speed of disk I/O.
38
39
40
41 326
42
43
44 327 Generally, there are two types of parallelism---data parallelism and task parallelism. The former
45
46
47
48 328 splits data horizontally into blocks of roughly equal sizes (the size of genomic intervals in our
49
50
51
52 329 case) before assigning them to all available processes; the latter assigns a roughly equal number of
53
54
55
56 330 input files to each process. For parallel multiway-merge, we choose data parallelism because the
57
58
59
60 331 implementation of task parallelism would be the same as the HPC-based implementation running
61
62
63
64
65

1 332 on a single node. Perhaps the most difficult part of data parallelism is uncertainty about the data
2
3
4
5 333 distribution across all input files, which usually leads to the problem of workload imbalance
6
7
8
9 334 among processes. If we pre-sample all the input files to estimate the record distribution, then a full
10
11
12
13 335 scan of the input files is required which will almost certainly takes more time than the single-
14
15
16
17 336 process multiway-merge method. As a compromise, we assume the distributions of SNP locations
18
19
20
21 337 in all VCF files are uniform and the input files can be split into regions of approximately equal
22
23
24
25 338 sizes. The total number of regions are set to be the number of concurrent processes, so that each
26
27
28
29 339 region is specifically handled by a process. To avoid seeking of a process's file reader to its
30
31
32
33 340 starting offset from the beginning of the file, we take advantage of the Tabix indexer [36], which
34
35
36
37 341 builds indices on data blocks of the input file and place the reader's pointer directly onto the
38
39
40
41 342 desired offset. One important aspect of the Tabix indexer is that it requires the input file to be
42
43
44
45 343 compressed in bgzip format which is not supported by Hadoop, HBase or Spark. The
46
47
48
49 344 compression and decompression of a file in bgzip format can be much faster than in bz2 format
50
51
52
53 345 used in our cluster-based schemas, single multiway-merge and HPC-based implementations, so
54
55
56
57 346 parallel multiway-merge can run much faster than other methods/schemas when input data size is
58
59
60 347 small.
61
62
63
64
65

1 348
2
3
4
5 349 For the HPC-based implementation, we adopt the task parallelism (Figure 5) to avoid sampling
6
7
8
9 350 and workload imbalance. Otherwise the workflow of HPC-based implementation is the same as
10
11
12
13 351 that of the MapReduce-based schema with the same operations and the same order: sampling in
14
15
16
17 352 parallel, dividing the dataset into splits of equal sizes, and assigning the splits to processes to
18
19
20
21 353 perform the merging. But this implementation is without data locality offered by HDFS and task
22
23
24
25 354 coordination offered by YARN and thus has a performance no better than the MapReduce-based
26
27
28
29 355 schema. Specifically, input files are shared across all nodes in the cluster via a Network File
30
31
32
33 356 System (NFS). In the first round, each core/process fetches roughly the same number of files from
34
35
36
37 357 the NFS and performs multiway-merging locally. In the following rounds, we adopted a tree-
38
39
40
41 358 structured execution strategy. In the second round, processes with even ID numbers (process id
42
43
44
45 359 starts from 0) retrieve the merged file from its adjacent process to the right, which are then merged
46
47
48
49 360 with its local merged file. Processes with odd ID number are terminated. In the third round,
50
51
52
53 361 processes with ID divisible by four retrieve the merged file from its adjacent process to the right in
54
55
56
57 362 the second round to merge with its local merged file. This process continues until all the files are
58
59
60
61 363 merged into a single file for a total of $\log(n)$ rounds, where n is the number of the input files.
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

364

365 **Strong and Weak Scalabilities**

366 In this study, we quantify scalability by measuring computing efficiency in tests of strong and

367 weak scalabilities. We define efficiency as the average time cost of processing a file per core:

368
$$\text{Efficiency} = (T_b * C_b / N_b) / (T_i * C_i / N_i)$$

369 where T_b is the baseline running time, C_b is the baseline number of cores, N_b is the baseline number

370 of input files, T_i is the current running time, C_i is the current number of cores, N_i is the current

371 number of input files. We also incorporated the parallel multiway-merge and MPI-based HPC

372 implementations as benchmarks in the tests.

373

374 For the strong scalability test, we fix the number of input files at 93 and increase the computing

375 resources up to 16-fold from the baseline. The baseline is a single node (4 cores) for all

376 methods/schemas except for the parallel multiway-merge in which only a single core is used

377 because it can only run on a single machine. For the weak scalability test, we increase both

378 computing resources and input data size at the same pace. The ratio is ten file/core for parallel

379 multiway-merge and ten file/node for all others.

1 380

2

3

4

5 381 **Results**

6

7

8

9 382 We conducted experiments of Apache cluster-based schemas using Amazon's Elastic MapReduce

10

11

12

13 383 (EMR) service and experiments of the HPC-based implementation using MIT's StarClusterTM

14

15

16

17 384 toolkit which launches an AWS openMP virtual private cluster (VPC). Within both infrastructures,

18

19

20

21 385 we choose EC2 working nodes of m3.xlarge type, which has four High Frequency Intel Xeon E5-

22

23

24

25 386 2670 v2 (Ivy Bridge) Processors and 15GB memory. We conducted experiments of parallel

26

27

28

29 387 multiway-merge on a single EC2 r4.8xlarge instance with 32 High Frequency Intel Xeon E5-2686

30

31

32

33 388 v4 (Broadwell) processors and 244 GB memory. We used a dataset consisting of 186 VCF files

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

390

391 **Overall Performance Analysis of Clustered-based Schemas**

392 Our primary goal is to explore the scalabilities of the three schemas on input data size and

393 available computing resources, namely CPUs. To achieve this, in this experiment we adjust the

394 number of input files from 10 to 186, with an approximate total uncompressed size from 2.5 GB to

395 40 GB, and used a varying number of working nodes from 3 to 18, namely 12 to 72 cores.

1 396
2
3
4
5 397 As Figure 6 shows, for all three schemas, given a fixed number of cores, the execution time
6
7
8
9 398 increases at a slower pace than that of the input data size. On the one hand, the increase of
10
11
12
13 399 execution time is more obvious with fewer cores because each core is fully utilized. As the
14
15
16
17 400 number of input files increases, so does the number of parallel tasks assigned to each core. For
18
19
20
21 401 example, given 12 cores, as the number of input files increases from 10 to 186 (18.6 fold), the
22
23
24
25 402 execution time increases from 739 to 4,366 seconds (~5.9 fold) for the MapReduce schema, from
26
27
28
29 403 375 to 5,479 seconds (~14.6 fold) for the HBase schema, and from 361 to 1,699 seconds (~4.7
30
31
32
33 404 fold) for the Spark schema. On the other hand, with relatively more cores such as 72, this linear
34
35
36
37 405 increasing trend is less pronounced because there are more cores than tasks so that all cores are
38
39
40
41 406 assigned at most one task. We also notice when input data size is small or moderate, the Spark
42
43
44
45 407 schema does not always show a consistent improvement in terms of execution time with more
46
47
48
49 408 cores. This is reflected, for example, in the intersection of curves occurred between 24 and 72
50
51
52
53 409 cores in Figure 6c. This phenomenon is attributed to the limitation of Spark's internal task
54
55
56
57 410 assignment policy which gives rise to the possibility that some nodes are assigned more than one
58
59
60
61 411 tasks while others remain idle.
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

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

Comparing Strong and Weak Scalabilities between Apache Cluster-based Schemas and Traditional Parallel Methods

Figure 7 shows the results of the strong scalability. In accordance with the Amdahl's law [37], all schemas/methods show degraded efficiency with increasing computing nodes/cores. Parallel multiway-merge has the steepest degradation because the more parallel processes, the higher likelihood of workload imbalances among them. In addition, disk I/O reaches saturation as more processes write simultaneously. Furthermore, to achieve data parallelism and improve execution speed, we used Tabix indexer to index data blocks of input files. While reading, each process needs to maintain a full copy of file descriptors, indices and uncompressed current data blocks of all input files in memory. When both the number of processes and input files are large, great pressure is placed on the memory management. For instance, a test with 93 files and 16 processes requires over 100GB memory, which results in a very long memory swap and garbage collection (GC) time. In contrast, the MapReduce-based schema has the best efficiency. Surprisingly, its efficiency even improves when the number of cores doubles from the baseline. This is because it has many parallel tasks in its second MapReduce phase, and when the core allowance is low, the

1 428 overheads of repetitive task launching and terminating on a single core become non-negligible.
2
3
4
5 429 Consequently, as the number of cores starts to increase, the actual proportion of overheads in the
6
7
8
9 430 total running time decreases, leading to an improved efficiency. Nonetheless, as the number of
10
11
12
13 431 cores further increases, the unparalleled parts of the schema gradually dominated the total running
14
15
16
17 432 time, leading to a reduced efficiency eventually.

18
19
20
21 433

22
23
24
25 434 For the weak scalability test (Figure 8), following Gustafson's law [38], all methods/schemas
26
27
28
29 435 show a much better efficiency than in the strong scalability test. Meanwhile, for the same reasons
30
31
32
33 436 as the strong scalability, the MapReduce-based schema enjoys the best efficiency while the HPC-
34
35
36
37 437 based implementation has the worst. This is because, for the HPC-based implementation, as the
38
39
40
41 438 number of input files increases, the total number of merging rounds also increases, leading to a
42
43
44
45 439 significantly reduced efficiency. Finally, all three Apache cluster-based schemas demonstrate
46
47
48
49 440 significantly better weak scalability than the two traditional parallel methods.

50
51
52 441

53 54 55 56 442 **The Anatomic Performances Analysis of Apache Cluster-based Schemas**

57
58
59
60 443 Another important goal of our study is to identify potential performance bottlenecks, so we
61
62
63
64
65

1 444 evaluate the execution time of each phase/stage of all three schemas. Figure 9 shows the trends of
2
3
4
5 445 the anatomic computing time spent on merging increasing number of VCF files (from 10 to 186)
6
7
8
9 446 using 48 cores. For the MapReduce schema (Figure 9a), its two phases account for a comparable
10
11
12
13 447 proportion of total time and both show a linear or sublinear scalability. The reason that the time
14
15
16
17 448 cost of the first phase between 40 and 93 input files remains flat is because both runs use two
18
19
20
21 449 rounds of mappers. As the number of files doubles to 186, four rounds of mappers are required
22
23
24
25 450 which results in about a two-fold increase in the time cost as expected. For the three phases of the
26
27
28
29 451 HBase schema (Figure 9b), they are scalable with input data size. Meanwhile, the second phase
30
31
32
33 452 becomes more dominant with more input files owing to the larger amount of shuffled data during
34
35
36
37 453 the writing of HFiles. However, we do not consider it as a bottleneck since all tasks of this phase
38
39
40
41 454 are parallelized with no workload or computational hotspot. We do not observe a super-linear
42
43
44
45 455 (relative to input data size) increment pattern from the figure neither. Finally, Figure 9c shows the
46
47
48
49 456 time costs of the three stages of the Spark schema. They show a uniform increasing trend with the
50
51
52
53 457 number of input files. Among them, the second stage takes up a considerable proportion of the
54
55
56
57 458 total execution time as it has a relatively expensive sort-by-key shuffling operation. Although no
58
59
60
61 459 data is shuffled in the first stage, its time lapse is close to the second stage. This is because at the
62
63
64
65

1 460 end of the first stage, data are sampled to determine the boundaries used by sort-by-key's range
2
3
4
5 461 partitioner. This operation demands a considerable execution time because it scans all the data and
6
7
8
9 462 balances them if necessary.
10
11
12
13 463
14
15
16
17 464 Given that no super-linear increasing trend is observed in running time for all phases/stages of the
18
19
20
21 465 three schemas, and they generally scale well with the input data size, we conclude although the
22
23
24
25 466 performances of these schemas might degrade to some extent when dealing with even larger input
26
27
28
29 467 data due to overheads such as data transmission over network, we would not expect any
30
31
32
33 468 significant bottleneck.

34
35
36
37 469

40 470 **Comparing Execution Speed between Apache Cluster-based Schemas and Traditional**

43 44 471 **Methods**

45
46
47
48 472 Another intriguing question is: how does the speed of the Apache cluster-based schemas compare
49
50
51
52 473 to single machine based and traditional parallel/distributed methods/applications on merging
53
54
55
56 474 multiple VCF files into a single VCF or TPED file? To answer this question, we choose the
57
58
59
60 475 widely-used VCFTools (v4.2) and a single-process multiway-merge implementation as single-

1 476 process benchmarks and parallel multiway-merge and HPC-based implementations as
2
3
4
5 477 parallel/distributed benchmarks, which are the same ones used in the experiments of strong and
6
7
8
9 478 weak scalabilities shown above.
10
11
12
13 479
14
15
16
17 480 In the first experiment, we merged 40 VCF files into one VCF file using VCFTools as the
18
19
20
21 481 benchmark. As shown in Table 2, VCFTools takes 30,189 seconds while the fastest Apache
22
23
24
25 482 cluster-based schema among the three, the MapReduce-based, takes only 484 seconds using 72
26
27
28
29 483 cores, representing about a 62-fold faster. In the second experiment (Figure 10), we tested the
30
31
32
33 484 time costs of merging of multiple VCF files into a single TPED file using single/parallel
34
35
36
37 485 multiway-merge and HPC-based implementations as benchmarks. The single multiway merger is
38
39
40
41 486 run on a node with the hardware configuration (4 cores and 15GB memory) identical to the nodes
42
43
44
45 487 on which the Apache cluster-based schemas are run. The parallel multiway merger is run on a
46
47
48
49 488 node with a maximum of 18 simultaneously running processes. The HPC-based implementation is
50
51
52
53 489 run on an 18-node cluster with the same hardware configuration as the cluster where the Apache
54
55
56
57 490 cluster-based schemas are run. Initially, with ten input files, the parallel multiway-merge (~30
58
59
60
61 491 seconds) is much faster than all the other methods: about 7.3-fold faster than the fastest Apache
62
63
64
65

1 492 cluster-based schema (MapReduce, 221 seconds). On the other hand, the slowest method is the
2
3
4
5 493 single-process multiway merger which takes 620 seconds to finish (about 2.8-fold slower than the
6
7
8
9 494 MapReduce-based schema). It is worth mentioning in this test the parallel multiway-merge is
10
11
12
13 495 essentially the same as the single-process multiway-merge, and the speed difference (~378
14
15
16
17 496 seconds) between them is the result of a different compression format (bz2 vs bgzip) of the input
18
19
20
21 497 files as explained above. As we gradually increase the number of input files to 186, the difference
22
23
24
25 498 in speed between the fastest overall method (parallel multiway merger, 602 seconds) and the
26
27
28
29 499 fastest Apache cluster-based schema (MapReduce, 809 seconds) reduces to about 1.3-fold, while
30
31
32
33 500 the difference between the slowest overall method (single multiway merger, 13,219 seconds) and
34
35
36
37 501 the MapReduce-based schema increases to 16.3-fold. In addition, all three Apache schemas
38
39
40
41 502 significantly outperform the HPC-based implementation. As explained in the strong and weak
42
43
44
45 503 scalabilities section above, we expect the larger the input data size, the faster the Apache cluster-
46
47
48
49 504 based schemas would run compared to the other traditional methods.
50
51
52
53 505
54
55
56 506 We also compare the time cost among the three schemas (Figure 10). They have a comparable
57
58
59
60 507 speed. More specifically, the MapReduce schema performs best if enough cores are available and
61
62
63
64
65

1 508 the input data size is large; the HBase schema performs best with moderate input data size; the
2
3
4
5 509 Spark schema performs best if only a limited number of cores are available and the input data size
6
7
8
9 510 is large. The rationale behind this observation is that, when the number of cores is sufficient, the
10
11
12
13 511 MapReduce-based schema can make the most use of the available computing resources because it
14
15
16
17 512 runs a constant 25 parallel jobs (one for each of chromosomes 1-22, X Y and M (Mitochondria))
18
19
20
21 513 in its second phase. In contrast, the Spark-based schema has fewer tasks whose number equals to
22
23
24
25 514 the number of input files to achieve maximum data-task locality. When the input data size is
26
27
28
29 515 moderate, the HBase-schema triumphs because its internal sorting and relative compact storage
30
31
32
33 516 format of intermediate data. When the input data size is large and computing resource is relatively
34
35
36
37 517 limited, the Spark-based schema outperforms the other two owing to its least number of data
38
39
40
41 518 shuffling (only one), execution plan optimization, and ability to cache intermediate results in
42
43
44
45 519 memory. We caution, however, the computing time may fluctuate depending on the distribution of
46
47
48
49 520 genomic locations in the input files as well as data loading balance of the HDFS.
50

51
52 521

53 54 55 56 522 **Discussion**

57
58
59
60 523 In this report, we describe three cluster-based schemas running on the Apache Hadoop
61
62
63
64
65

1 524 (MapReduce), HBase and Spark platforms respectively for performing sorted merging of variants
2
3
4
5 525 identified from WGS. We show all three schemas are scalable on both input data size and
6
7
8
9 526 computing resources, suggesting large scale of ‘-omics’ data can be merged efficiently given the
10
11
12
13 527 computing resources readily available in the cloud. Furthermore, the three schemas show better
14
15
16
17 528 strong and weak scalabilities than traditional single machine-based parallel multiway-merge and
18
19
20
21 529 cluster-based HPC methods owing to the absence of I/O bottleneck, better workload balance
22
23
24
25 530 among nodes, less pressure on memory, as well as data locality and efficient task coordination
26
27
28
29 531 mechanisms provided by HDFS and YARN. We also show even with a moderate-sized cluster and
30
31
32
33 532 input data, all three schemas significantly outperform the broadly-used, single-machine based
34
35
36
37 533 VCFTools, single-process multiway-merge and HPC-based implementations. Although initially
38
39
40
41 534 the parallel multiway-merge implementation is much faster than the Apache schemas owing to its
42
43
44
45 535 advantage of local I/O and light compression of input files, its poor scalability diminishes its
46
47
48
49 536 initial advantage as the number of concurrent processes and input files increases. Consequently,
50
51
52
53 537 we expect the Apache cluster-based schemas eventually outperform the parallel multiway-merge
54
55
56
57 538 when merging a much larger scale of data using a larger number of cores.
58
59
60
61
62
63
64
65

1 539
2
3
4
5
6 540 Unlike normal merging, efficient sorted merging of many large tables has always been a difficult
7
8
9
10 541 problem in the field of data management. Multiway-merge is the most efficient single-machine
11
12
13 542 based method for sorted merging, but its performance is limited by the disk I/O [39]. Sorted
14
15
16
17 543 merging also places challenges to distributed system based solutions because neither the efficient
18
19
20
21 544 hash-based merging nor caching the intermediate table in shared memory is feasible [40].
22
23
24
25 545 Although a utility named *total-order-joining* is provided by the Hadoop for addressing this
26
27
28
29 546 problem, it suffers from both network communication and local disk I/O bottlenecks, thus is not
30
31
32
33 547 scalable [27, 41]. In contrast, our schemas divide this problem into different phases/stages of tasks
34
35
36
37 548 each conquered in parallel to bypass these bottlenecks and achieve maximum parallelism.
38
39
40
41 549 Furthermore, in addition to merging sequencing variant data, these schemas can be generalized for
42
43
44
45 550 other key-based, sorted merging problems are frequently encountered in genetics and genomics
46
47
48
49 551 data processing. As an example, they can be slightly modified to merge multiple BED format files
50
51
52
53 552 such as ChIP-seq peak lists [42] and other genomic regions of interest. Other potentially useful
54
55
56
57 553 features include: 1) Unlike traditional sorted merging algorithms which usually require presorted
58
59
60
61
62
63
64
65

1 554 inputs for a better performance, our schemas are free of such a requirement; 2) Our
2
3
4
5 555 implementations automatically take care of multi-allelic positions which are frequent in large scale
6
7
8
9 556 VCF files by retaining the information of all alleles until the merging actually occurs.
10
11
12
13 557
14
15
16
17 558 Finally, in light of these different features and specialties of these three platforms, each of the
18
19
20
21 559 three schemas we developed has its own advantages and disadvantages under different application
22
23
24
25 560 scenarios as summarized in Table 1. For example, the MapReduce schema is good for a static one-
26
27
28
29 561 time, non-incremental merging on large-size data provided sufficient cores are available since it
30
31
32
33 562 has the most parallel jobs, the least overheads, and the most transparent workflow. The HBase
34
35
36
37 563 schema, supported by data warehousing technologies, fits for an incremental merging since it does
38
39
40
41 564 not need to re-merge existing results with new ones from the scratch only if the incremental
42
43
44
45 565 merging is performed on the same chromosomes. Also, it provides a highly-efficient storage and
46
47
48
49 566 On-Line Analytical Processing (OLAP) on merged results. The Spark schema is ideal for merging
50
51
52
53 567 large scale data with relatively limited computing resources because it has the least data shuffling
54
55
56
57 568 and keeps intermediate results in memory. A bonus brought by Spark is the subsequent statistical
58
59
60
61 569 analyses can be carried out directly on the merged results using its rich set of parallel statistical
62
63
64
65

1 570 utilities.
2
3

4
5
6 571
7

8
9
10 572 **Availability and Requirements**
11

12
13 573 Project name: CloudMerge
14
15

16
17 574 Project home page: <https://github.com/xsun28/CloudMerge>
18
19

20
21 575 Operating system(s): Linux
22
23

24
25 576 Programming language: Java, Python
26
27

28
29 577 Other requirements: Java 1.7 or higher, Python 2.7 or 3.6, Hadoop-2.7, HBase-1.3, Spark-2.1,
30
31

32
33 578 StarCluster 0.95, MPI for Python 3.0.0
34
35

36
37 579 License: Apache License 2.0
38
39

40
41 580
42
43

44
45 581 **Availability of Data and Materials**
46
47

48
49 582 The source codes of the project are available in GitHub. The 186 individual VCF files used in our
50
51

52
53 583 experiments are modified from the original VCF files obtained from WGS conducted by the
54
55

56
57 584 Consortium on Asthma among African-ancestry Population in the Americas (CAAPA) [32]. To
58
59
60
61
62
63
64
65

1 585 conceal the potential individual identifiable genotype information from the public, we encrypt the
2
3
4
5 586 authentic genomic location of the original 93 VCF files to generate a new batch of encrypted VCF
6
7
8
9 587 files for test purposes. Please refer to *Data Formats and Operations* section for details. These
10
11
12
13 588 supporting data and a snapshot of project codes are available at the *GigaScience* database,
14
15
16
17 589 GigaDB [43]. Via GigaDB, we also provide sample results of merging 93VCF files into either one
18
19
20
21 590 VCF or one TPED file using our Apache cluster-based schemas.
22
23
24

25 591

28 592 **Abbreviations**

30
31
32 593 VCF: Variant Call Format; MPI: Message Passing Interface; HPC: High Performance Computing;
33
34
35
36 594 GWAS: Genome Wide Association Studies; WGS: Whole Genome Sequencing; WES: whole
37
38
39
40 595 exome sequencing; AWS: Amazon Web Service; CGDM: Collaborative Genomic Data Model;
41
42
43
44 596 SAM/BAM: Sequence/Binary Alignment/Map; RDD: Resilient Distributed Dataset; DAG: Directed
45
46
47
48 597 Acyclic Graph; SPMD: Single Program, Multiple Data; NFS: Network File System; EMR: Elastic-
49
50
51
52 598 MapReduce; VPC: Virtual Private Cluster; GC: Garbage Collection; CAAPA: Consortium on
53
54
55
56 599 Asthma among African-ancestry Population in the Americas;
57
58
59

60 600
61
62
63
64
65

1 601 **Ethics Approval and Consent to Participate**

2
3
4
5 602 Ethics approval for the CAAPA program was provided by the Johns Hopkins University
6
7
8
9 603 Institutional Review Board following commencement of the study in 2011 (IRB00045892 The
10
11
12
13 604 Consortium on Asthma among African-ancestry Populations), and included study team members
14
15
16
17 605 from each CAAPA site, including Emory University (site PI, Zhaohui Qin). Access to the raw data
18
19
20
21 606 as CAAPA team members is granted according to the guideline of the IRB-approved study. Informed
22
23
24
25 607 consent has been obtained from all study participants of CAAPA.
26
27

28 608 **Competing Interests**

29 609 The authors declare they have no competing interests.
30
31
32
33
34
35
36
37 610

38
39
40 611 **Funding**

41
42
43
44 612 This study was supported by grants from National Heart, Lung, and Blood Institute [R01HL104608,
45
46
47
48 613 R01HL117004, R01HL128439, R01HL135156, X01HL134589]; National Institute of
49
50
51
52 614 Environmental Health Sciences [R01ES015794, R21ES24844]; National Institute on Minority
53
54
55
56 615 Health and Health Disparities [P60MD006902, R01MD010443, RL5GM118984]; National
57
58
59
60 616 Institute of Neurological Disorders and Stroke [R01NS051630, P01NS097206, U54NS091859];
61
62
63
64
65

1 617 National Science Foundation [ACI 1443054, IIS 1350885,]; Tobacco-Related Disease Research
2
3
4
5 618 Program [24RT-0025]. The Genes-Environments and Admixture in Latino Americans (GALA II)
6
7
8
9 619 Study, the Study of African Americans, Asthma, Genes and Environments (SAGE) Study and E.G.B.
10
11
12
13 620 are supported by the Sandler Family Foundation, the American Asthma Foundation, the RWJF
14
15
16
17 621 Amos Medical Faculty Development Program and the Harry Wm. and Diana V. Hind Distinguished
18
19
20
21 622 Professor in Pharmaceutical Sciences II.
22
23
24 623

26 624 **Authors' Contributions**

28
29
30 625 J.G. introduced the problem. X.S., F.W. initiated this project. X.S. designed and implemented the
31
32
33
34 626 CloudMerge project. X.S. drafted the manuscript. X.S., J.P., F.W. and Z.Q. revised the manuscript.
35
36
37
38 627 K.C.B. conceived the initial consortium design, acquired biospecimens for NGS, facilitated generation
39
40
41
42 628 of NGS data. K.C.B., R.A.M., I.R., T.H.B. conceived initial experiments, interpreted NGS data.
43
44
45
46 629 E.G.B., C.E. acquired biospecimens for NGS, facilitated generation of NGS data.
47
48

49
50 630

54 631 **Acknowledgements**

55
56
57
58 632 We thank the three referees for their constructive critiques and detailed comments. We are grateful to Ms.
59
60
61
62
63
64
65

1 633 Mary Taylor Mann and Ms. Alyssa Leann Duck for their editorial help during writing and revising of the
2
3
4
5 634 manuscript. E.G.B. wish to acknowledge the following GALA II and SAGE co-investigators for subject
6
7
8
9 635 recruitment, sample processing and quality control: Sandra Salazar, Scott Huntsman, MSc, Donglei Hu,
10
11
12
13 636 PhD, Lisa Caine, Shannon Thyne, MD, Harold J. Farber, MD, MSPH, Pedro C. Avila, MD, Denise
14
15
16
17 637 Serebrisky, MD, William Rodriguez-Cintron, MD, Jose R. Rodriguez-Santana, MD, Rajesh Kumar, MD,
18
19
20
21 638 Luisa N. Borrell, DDS, PhD, Emerita Brigino-Buenaventura, MD, Adam Davis, MA, MPH, Michael A.
22
23
24
25 639 LeNoir, MD, Kelley Meade, MD, Saunak Sen, PhD and Fred Lurmann, MS.
26
27
28
29 640
30
31
32

33 641 **References**

- 34
35
36 642 1. Massie M, Nothaft F, Hartl C, Kozanitis C, Schumacher A, Joseph AD, et al. Adam: Genomics formats
37
38 643 and processing patterns for cloud scale computing. University of California, Berkeley Technical
39
40 644 Report, No UCB/Eecs-2013. 2013;207.
- 41
42 645 2. Siretskiy A, Sundqvist T, Voznesenskiy M and Spjuth O. A quantitative assessment of the hadoop
43
44 646 framework for analyzing massively parallel dna sequencing data. Gigascience. 2015;4 1:26.
- 45
46 647 3. Merelli I, Pérez-Sánchez H, Gesing S and D'Agostino D. Managing, analysing, and integrating big
47
48 648 data in medical bioinformatics: open problems and future perspectives. BioMed research
49
50 649 international. 2014;2014.
- 51
52 650 4. Reyes-Ortiz JL, Oneto L and Anguita D. Big data analytics in the cloud: Spark on hadoop vs
53
54 651 mpi/openmp on beowulf. Procedia Computer Science. 2015;53:121-30.
- 55
56 652 5. Burren OS, Guo H and Wallace C. VSEAMS: a pipeline for variant set enrichment analysis using
57
58 653 summary GWAS data identifies IKZF3, BATF and ESRRA as key transcription factors in type 1
59
60 654 diabetes. Bioinformatics. 2014;30 23:3342-8. doi:10.1093/bioinformatics/btu571.
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
- 655 6. Apache Hadoop. <http://hadoop.apache.org/>. Accessed 10 Oct 2017.
 - 656 7. Dean J and Ghemawat S. Mapreduce: Simplified data processing on large clusters. *Commun Acm.*
657 2008;51 1:107-13. doi:Doi 10.1145/1327452.1327492.
 - 658 8. Vora MN. Hadoop-HBase for large-scale data. In: *Computer science and network technology*
659 *(ICCSNT), 2011 international conference on 2011*, pp.601-5. IEEE.
 - 660 9. Zaharia M, Chowdhury M, Das T, Dave A, Ma J, McCauley M, et al. Resilient distributed datasets: A
661 fault-tolerant abstraction for in-memory cluster computing. In: *Proceedings of the 9th USENIX*
662 *conference on Networked Systems Design and Implementation 2012*, pp.2-. USENIX Association.
 - 663 10. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernysky A, et al. The Genome Analysis
664 Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome*
665 *research.* 2010;20 9:1297-303.
 - 666 11. Schatz MC. CloudBurst: highly sensitive read mapping with MapReduce. *Bioinformatics.* 2009;25
667 11:1363-9. doi:10.1093/bioinformatics/btp236.
 - 668 12. Pireddu L, Leo S and Zanetti G. SEAL: a distributed short read mapping and duplicate removal tool.
669 *Bioinformatics.* 2011;27 15:2159-60.
 - 670 13. Niemenmaa M, Kallio A, Schumacher A, Klemelä P, Korpelainen E and Heljanko K. Hadoop-BAM:
671 directly manipulating next generation sequencing data in the cloud. *Bioinformatics.* 2012;28
672 6:876-7.
 - 673 14. Langmead B, Schatz MC, Lin J, Pop M and Salzberg SL. Searching for SNPs with cloud computing.
674 *Genome Biol.* 2009;10 11:R134. doi:10.1186/gb-2009-10-11-r134.
 - 675 15. Wang S, Mares MA and Guo YK. CGDM: collaborative genomic data model for molecular profiling
676 data using NoSQL. *Bioinformatics.* 2016;32 23:3654-60. doi:10.1093/bioinformatics/btw531.
 - 677 16. Zhou W, Li R, Yuan S, Liu C, Yao S, Luo J, et al. MetaSpark: a spark-based distributed processing tool
678 to recruit metagenomic reads to reference genomes. *Bioinformatics.* 2017;33 7:1090-2.
 - 679 17. Niu B, Zhu Z, Fu L, Wu S and Li W. FR-HIT, a very fast program to recruit metagenomic reads to
680 homologous reference genomes. *Bioinformatics.* 2011;27 12:1704-5.
 - 681 18. AWS Genomics Guide.
682 https://d0.awsstatic.com/Industries/HCLS/Resources/AWS_Genomics_WP.pdf. Accessed
683 10 Oct 2017.
 - 684 19. Gruber K. Google for genomes. *Nature Research*, 2014.
 - 685 20. O'Brien AR, Saunders NF, Guo Y, Buske FA, Scott RJ and Bauer DC. VariantSpark: population scale

- 686 clustering of genotype information. BMC genomics. 2015;16 1:1052.
- 687 21. Wiewiórka MS, Messina A, Pacholewska A, Maffioletti S, Gawrysiak P and Okoniewski MJ. SparkSeq:
688 fast, scalable and cloud-ready tool for the interactive genomic data analysis with nucleotide
689 precision. Bioinformatics. 2014;30 18:2652-3.
- 690 22. Decap D, Reumers J, Herzeel C, Costanza P and Fostier J. Halvade: scalable sequence analysis with
691 MapReduce. Bioinformatics. 2015;31 15:2482-8.
- 692 23. He M, Person TN, Hebbring SJ, Heinzen E, Ye Z, Schrodi SJ, et al. SeqHBase: a big data toolset for
693 family based sequencing data analysis. Journal of medical genetics. 2015:jmedgenet-2014-
694 102907.
- 695 24. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, et al. PLINK: a tool set for
696 whole-genome association and population-based linkage analyses. Am J Hum Genet. 2007;81
697 3:559-75. doi:10.1086/519795.
- 698 25. Mohammed EA, Far BH and Naugler C. Applications of the MapReduce programming framework
699 to clinical big data analysis: current landscape and future trends. BioData Min. 2014;7:22.
700 doi:10.1186/1756-0381-7-22.
- 701 26. Huang H, Tata S and Prill RJ. BlueSNP: R package for highly scalable genome-wide association
702 studies using Hadoop clusters. Bioinformatics. 2013;29 1:135-6.
703 doi:10.1093/bioinformatics/bts647.
- 704 27. White T. Hadoop: The definitive guide. " O'Reilly Media, Inc."; 2012.
- 705 28. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al. The variant call format and
706 VCFtools. Bioinformatics. 2011;27 15:2156-8. doi:10.1093/bioinformatics/btr330.
- 707 29. Multiway-Merge Algorithm. https://en.wikipedia.org/wiki/K-Way_Merge_Algorithms.
708 Accessed 10 Oct 2017.
- 709 30. Chang F, Dean J, Ghemawat S, Hsieh WC, Wallach DA, Burrows M, et al. Bigtable: A distributed
710 storage system for structured data. Acm T Comput Syst. 2008;26 2 doi:Artn 4
711 10.1145/1365815.1365816.
- 712 31. Genomes Project C, Abecasis GR, Altshuler D, Auton A, Brooks LD, Durbin RM, et al. A map of
713 human genome variation from population-scale sequencing. Nature. 2010;467 7319:1061-73.
714 doi:10.1038/nature09534.
- 715 32. Mathias RA, Taub MA, Gignoux CR, Fu W, Musharoff S, O'Connor TD, et al. A continuum of
716 admixture in the Western Hemisphere revealed by the African Diaspora genome. Nat Commun.

1 717 2016;7:12522. doi:10.1038/ncomms12522.

2 718 33. Kwon Y, Balazinska M, Howe B and Rolia J. A study of skew in mapreduce applications. Open Cirrus

3

4 719 Summit. 2011;11.

5

6 720 34. O'Neil P, Cheng E, Gawlick D and O'Neil E. The log-structured merge-tree (LSM-tree). Acta

7

8 721 Informatica. 1996;33 4:351-85.

9

10 722 35. CloudMerge. <https://github.com/xsun28/CloudMerge>

11

12 723 36. Li H. Tabix: fast retrieval of sequence features from generic TAB-delimited files. Bioinformatics.

13

14 724 2011;27 5:718-9.

15

16 725 37. Amdahl GM. Validity of the single processor approach to achieving large scale computing

17

18 726 capabilities. In: *Proceedings of the April 18-20, 1967, spring joint computer conference 1967*,

19

20 727 pp.483-5. ACM.

21

22 728 38. Gustafson JL. Reevaluating Amdahl's law. *Commun Acm*. 1988;31 5:532-3.

23

24 729 39. Sedgewick R and Flajolet P. An introduction to the analysis of algorithms. Addison-Wesley; 2013.

25

26 730 40. Özsü MT and Valduriez P. Principles of distributed database systems. Springer Science & Business

27

28 731 Media; 2011.

29

30 732 41. Miner D and Shook A. MapReduce Design Patterns: Building Effective Algorithms and Analytics for

31

32 733 Hadoop and Other Systems. " O'Reilly Media, Inc."; 2012.

33

34 734 42. Chen L, Wang C, Qin ZS and Wu H. A novel statistical method for quantitative comparison of

35

36 735 multiple ChIP-seq datasets. *Bioinformatics*. 2015;31 12:1889-96.

37

38 736 43. Sun X, Gao J, Jin P, Celeste Eng, Esteban G. Burchard, Terri H. Beaty, Ingo Ruczinski, Rasika A.

39

40 737 Mathias, Kathleen C. Barnes, Wang F and Qin Z. Supporting data for "Optimized Distributed

41

42 738 Systems Achieve Significant Performance Improvement on Sorted Merging of Massive VCF

43

44 739 Files" GigaScience Database 2018. <http://dx.doi.org/10.5524/100423>.

45

46 740

47

48

49 741

50

51

52 742

53

54

55

56 743

57

58

59

60 744

61

62

63

64

65

1 **745 Complete List of CAAPA Consortium Members and Their Affiliations**

- 2
3
4
5 746 Kathleen C. Barnes, PhD^{1:2}, Terri H. Beaty, PhD², Meher Preethi Boorgula MS¹, Monica
6
7
8
9 747 Campbell, BS¹, Sameer Chavan, MS¹, Jean G. Ford, MD^{2:3}, Cassandra Foster, CCRP¹, Li Gao,
10
11
12 748 MD, PhD¹, Nadia N. Hansel, MD, MPH¹, Edward Horowitz, BA¹, Lili Huang, MPH¹, Rasika Ann
13
14
15 749 Mathias, ScD^{1:2}, Romina Ortiz, MA¹, Joseph Potee, MS¹, Nicholas Rafaels, MS¹, Ingo Ruczin-
16
17
18
19
20
21 750 ski, PhD⁴, Alan F. Scott, PhD¹, Margaret A. Taub, PhD⁴, Candelaria Ver-gara, PhD¹, Jingjing
22
23
24
25 751 Gao, PhD⁵, Yijuan Hu, PhD⁶, Henry Richard Johnston, PhD⁶, Zhaohui S. Qin, PhD⁶, Albert M.
26
27
28
29 752 Levin, PhD⁷, Badri Padhukasahas-ram, PhD⁸, L. Keoki Williams, MD, MPH^{8:9}, Georgia M.
30
31
32
33 753 Dunston, PhD^{10:11}, Mezbah U. Faruque, MD, PhD¹¹, Eimear E. Kenny, PhD^{12:13}, Kimberly Gi-
34
35
36
37 754 etzen, PhD¹⁴, Mark Hansen, PhD¹⁴, Rob Genuario, PhD¹⁴, Dave Bullis, MBA¹⁴, Cindy Lawley,
38
39
40
41 755 PhD¹⁴, Aniket Deshpande, MS¹⁵, Wendy E. Grus, PhD¹⁵, Devin P. Locke, PhD¹⁵, Marilyn G.
42
43
44
45 756 Foreman, MD¹⁶, Pedro C. Avila, MD¹⁷, Leslie Grammer, MD¹⁷, Kwang-Youn A. Kim, PhD¹⁸,
46
47
48
49 757 Rajesh Kumar, MD^{19:20}, Robert Schleimer, PhD²¹, Carlos Bustamante, PhD¹², Francisco
50
51
52
53 758 M. De La Vega, DS¹², Chris R. Gignoux, MS¹², Suyash S. Shringarpure, PhD¹², Shaila Musharo,
54
55
56
57 759 MS¹², Genevieve Wojcik, PhD¹², Esteban G. Burchard, MD, MPH^{22:23}, Celeste Eng, BS²³, Pierre-
58
59
60
61 760 Antoine Gourraud, PhD²⁴, Ryan D. Hernandez, PhD^{22:25:26}, Antoine Lizee, PhD²⁴, Maria Pino-

- 1 761 Yanes, PhD^{23:27}, Dara G. Torgerson, PhD²³, Zachary A. Szpiech, PhD²², Raul Torres, BS²⁸, Dan L.
2
3
4
5 762 Nicolae, PhD^{29:30}, Carole Ober, PhD³¹, Christopher O Olopade, MD, MPH³², Olufunmilayo
6
7
8
9 763 Olopade, MD²⁹, Oluwafemi Oluwole, MSc²⁹, Ganiyu Arinola, PhD³³, Timothy D. O'Connor,
10
11
12
13 764 PhD^{34:35:36}, Wei Song, PhD^{34:35:36}, Goncalo Abecasis, DPhil³⁷, Adolfo Correa, MD, MPH, PhD³⁸,
14
15
16
17 765 Solomon Musani, PhD³⁸, James G. Wilson, MD³⁹, Leslie A. Lange, PhD⁴⁰, Joshua Akey, PhD⁴¹,
18
19
20
21 766 Michael Bamshad, MD⁴², Jessica Chong, PhD⁴², Wenqing Fu, PhD⁴¹, Deborah Nickerson, PhD⁴¹,
22
23
24
25 767 Alexander Reiner, MD, MSc⁴³, Tina Hartert, MD, MPH⁴⁴, Lorraine B. Ware, MD^{44:45}, Eugene
26
27
28
29 768 Bleecker, MD⁴⁶, Deborah Meyers, PhD⁴⁶, Victor E. Ortega, MD⁴⁶, Pissamai Maul, BSc, RN⁴⁷,
30
31
32
33 769 Trevor Maul, RN⁴⁷, Harold Watson, MD^{48:49}, Maria Ilma Araujo, MD, PhD⁵⁰, Ricardo Riccio
34
35
36
37 770 Oliveira, PhD⁵¹, Luis Caraballo, MD, PhD⁵², Javier Marrugo, MD⁵³, Beatriz Martinez, MSc⁵²,
38
39
40
41 771 Catherine Meza, LB⁵², Gerardo Ayestas⁵⁴, Edwin Francisco Herrera-Paz, MD, MSc^{55:56:57}, Pamela
42
43
44
45 772 Landaverde-Torres⁵⁵, Said Omar Leiva Erazo⁵⁵, Rosella Martinez, BSc⁵⁵, varo Mayorga, MD⁵⁶,
46
47
48
49 773 Luis F. Mayorga, MD⁵⁵, Delmy-Aracely Mejia-Mejia, MD^{56:57}, Hector Ramos⁵⁵, Allan Saenz⁵⁴,
50
51
52
53 774 Gloria Varela⁵⁴, Olga Marina Vasquez⁵⁷, Trevor Ferguson, MBBS, DM, MSc⁵⁸, Jennifer Knight-
54
55
56
57 775 Madden, MBBS, PhD⁵⁸, Maureen Samms-Vaughan, MBBS, DM, Ph⁵⁹, Rainford J. Wilks, MBBS,
58
59
60
61
62
63
64
65

- 1 776 DM, MSc⁵⁸, Akim Adegnika, MD, PhD^{60;61;62}, Ulysse Ateba-Ngoa, MD^{60;61;62}, Maria
2
3
4
5 777 Yazdanbakhsh, PhD⁶²
6
7
8
9 778 1 Department of Medicine, Johns Hopkins University, Baltimore, MD.
10
11
12
13 779 2 Department of Epidemiology, Bloomberg School of Public Health, JHU, Baltimore, MD.
14
15
16
17 780 3 Department of Medicine, The Brooklyn Hospital Center, Brooklyn, NY.
18
19
20
21 781 4 Department of Biostatistics, Bloomberg School of Public Health, JHU, Baltimore, MD.
22
23
24
25 782 5 Data and Statistical Sciences, AbbVie, North Chicago, IL.
26
27
28
29 783 6 Department of Biostatistics and Bioinformatics, Emory University, Atlanta, GA.
30
31
32
33 784 7 Department of Public Health Sciences, Henry Ford Health System, Detroit, MI.
34
35
36
37 785 8 Center for Health Policy & Health Services Research, Henry Ford Health System, Detroit, MI.
38
39
40
41 786 9 Department of Internal Medicine, Henry Ford Health System, Detroit, MI.
42
43
44
45 787 10 Department of Microbiology, Howard University College of Medicine, Washington, DC.
46
47
48
49 788 11 National Human Genome Center, Howard University College of Medicine, Washington, DC.
50
51
52
53 789 12 Department of Genetics, Stanford University School of Medicine, Stanford, CA.
54
55
56
57 790 13 Department of Genetics and Genomics, Icahn School of Medicine at Mount Sinai, New York.
58
59
60
61 791 14 Illumina, Inc., San Diego, CA.
62
63
64
65

- 1 792 15 Knome Inc., Cambridge, MA.
2
3
4
5 793 16 Pulmonary and Critical Care Medicine, Morehouse School of Medicine, Atlanta, GA.
6
7
8
9 794 17 Department of Medicine, Northwestern University, Chicago, IL.
10
11
12
13 795 18 Department of Preventive Medicine, Northwestern University, Chicago, IL.
14
15
16
17 796 19 Department of Pediatrics, Northwestern University, Chicago, IL.
18
19
20
21 797 20 The Ann & Robert H. Lurie Children's Hospital of Chicago, Chicago.
22
23
24
25 798 21 Department of Medicine, Northwestern Feinberg School of Medicine, Chicago,
26
27
28
29 799 22 Department of Bioengineering and Therapeutic Sciences, University of California, San
30
31
32
33 800 Francisco, San Francisco, CA.
34
35
36
37 801 23 Department of Medicine, University of California, San Francisco, San Francisco, CA.
38
39
40
41 802 24 Department of Neurology, University of California, San Francisco, San Francisco, CA.
42
43
44
45 803 25 Institute for Human Genetics, Institute for Human Genetics, University of California, San
46
47
48
49 804 Francisco, San Francisco.
50
51
52
53 805 26 California Institute for Quantitative Biosciences, University of California, San Francisco , San
54
55
56
57 806 Francisco, CA.
58
59
60
61 807 27 CIBER de Enfermedades Respiratorias, Instituto de Salud Carlos III, Madrid.
62
63
64
65

- 1 808 28 Biomedical Sciences Graduate Program, University of California, San Francisco, San
2
3
4
5 809 Francisco, CA.
6
7
8
9 810 29 Department of Medicine, University of Chicago, Chicago, IL.
10
11
12
13 811 30 Department of Statistics, University of Chicago, Chicago, IL.
14
15
16
17 812 31 Department of Human Genetics, University of Chicago, Chicago, IL.
18
19
20
21 813 32 Department of Medicine and Center for Global Health, University of Chicago, Chicago, IL.
22
23
24
25 814 33 Department of Chemical Pathology, University of Ibadan, Ibadan, Nigeria.
26
27
28
29 815 34 Institute for Genome Sciences, University of Maryland School of Medicine, Baltimore, MD.
30
31
32
33 816 35 Program in Personalized and Genomic Medicine, University of Maryland School of Medicine,
34
35
36
37 817 Baltimore.
38
39
40
41 818 36 Department of Medicine, University of Maryland School of Medicine, Baltimore, MD.
42
43
44
45 819 37 Department of Biostatistics, SPH II, University of Michigan, Ann Arbor, MI.
46
47
48
49 820 38 Department of Medicine, University of Mississippi Medical Center, Jackson, MS.
50
51
52
53 821 39 Department of Physiology and Biophysics, University of Mississippi Medical Center, Jackson,
54
55
56
57 822 MS.
58
59
60 823 40 Department of Genetics, University of North Carolina, Chapel Hill, NC.
61
62
63
64
65

- 1 824 41 Department of Genomic Sciences, University of Washington, Seattle, WA.
2
3
4
5 825 42 Department of Pediatrics, University of Washington, Seattle, WA.
6
7
8
9 826 43 University of Washington, Seattle, WA.
10
11
12
13 827 44 Department of Medicine, Vanderbilt University, Nashville, TN.
14
15
16
17 828 45 Department of Pathology, Microbiology and Immunology, Vanderbilt University, Nashville.
18
19
20
21 829 46 Center for Human Genomics and Personalized Medicine, Wake Forest School of Medicine,
22
23
24
25 830 Winston-Salem, NC.
26
27
28
29 831 47 Genetics and Epidemiology of Asthma in Barbados, The University of the West Indies.
30
31
32
33 832 48 Faculty of Medical Sciences Cave Hill Campus, The University of the West Indies.
34
35
36
37 833 49 Queen Elizabeth Hospital, Queen Elizabeth Hospital, The University of the West Indies.
38
39
40
41 834 50 Immunology Service, Universidade Federal da Bahia, Salvador, BA.
42
43
44
45 835 51 Laboratrio de Patologia Experimental, Centro de Pesquisas Gonalo Moniz, Salvador, BA.
46
47
48
49 836 52 Institute for Immunological Research, Universidad de Cartagena, Cartagena.
50
51
52
53 837 53 Instituto de Investigaciones Immunologicas, Universidad de Cartagena, Cartagena.
54
55
56
57 838 54 Faculty of Medicine, Universidad Nacional Autonoma de Honduras en el Valle de Sula, San
58
59
60 839 Pedro Sula.
61
62
63
64
65

1 840 55 Facultad de Medicina, Universidad Catolica de Honduras, San Pedro Sula.

2
3
4
5 841 56 Centro de Neumologia y Alergias, San Pedro Sula.

6
7
8
9 842 57 Faculty of Medicine, Centro Medico de la Familia, San Pedro Sula.

10
11
12
13 843 58 Tropical Medicine Research Institute, The University of the West Indies.

14
15
16
17 844 59 Department of Child Health, The University of the West Indies.

18
19
20
21 845 60 Centre de Recherches Mdicales de Lambarn.

22
23
24
25 846 61 Institut fr Tropenmedizin, Universitt Tbingen.

26
27
28
29 847 62 Department of Parasitology, Leiden University Medical Center, Netherlands.

30
31
32
33 848

34
35
36
37 849 **Figure legends**

38
39
40
41 850 **Figure 1. Merging multiple VCF files into a single TPED file.** Left tables represent input VCF

42
43
44 851 files. Table to the right represents the merged TPED file. Records are filtered out if their Filter

45
46
47
48 852 value does not equal to "PASS" (Pos 10147). Individual genotypes from multiple VCF files with

49
50
51
52 853 the same genomic location are aggregated together in one row. The resulting TPED file thus has

53
54
55
56 854 an inclusive set of sorted genomic locations of all variants found in the input VCF files.

57
58
59
60 855

61
62
63
64
65

1 856 **Figure 2. The workflow chart of the MapReduce schema.** The workflow is divided into two
2
3
4
5 857 phases: In the first phase, variants are filtered, grouped by chromosomes into bins, and mapped
6
7
8
9 858 into key-value records. Two sampling steps are implemented to generate partition lists of all
10
11
12
13 859 chromosomes. In the second phase, parallel jobs of specified chromosomes are launched. Within
14
15
16
17 860 each job, records from corresponding bins are loaded, partitioned, sorted and merged by genomic
18
19
20
21 861 locations before being saved into a TPED file.

22
23
24
25 862

26
27
28
29 863 **Figure 3. The workflow chart of the HBase schema.** The workflow is divided into three phases.
30
31
32
33 864 The first is a sampling, filtering and mapping phase. A MapReduce job samples out variants
34
35
36
37 865 whose genomic positions are used as region boundaries when creating the HBase table. Only
38
39
40
41 866 qualified records are mapped as key-values and saved as Hadoop sequence files. The second is the
42
43
44
45 867 HBase bulk loading phase in which a MapReduce job loads and writes records generated from the
46
47
48
49 868 previous phase, aggregating them into corresponding regional HFiles in the form of HBase's row
50
51
52
53 869 key and column families. Finished HFiles are moved into HBase data storage folders on region
54
55
56
57 870 servers. In the third phase, parallel scans were launched over regions of the whole table to retrieve
58
59
60
61 871 desired records which are subsequently merged and exported to the TPED file.

1 872
2
3
4
5 873 **Figure 4. The workflow chart of the Spark schema.** The workflow is divided into three stages.
6
7
8
9 874 In the first stage, VCF records are loaded, filtered, and mapped to PairRDDs with keys of genomic
10
11
12
13 875 position and values of genotype. The sort-by-key shuffling spans across the first two stages,
14
15
16
17 876 sorting and grouping together records by keys. Then grouped records with the same key are
18
19
20
21 877 locally merged into one record in TPED format. Finally, merged records are exported to the TPED
22
23
24
25 878 file.

26
27
28 879
29
30
31
32 880 **Figure 5. The execution plan of the HPC-based implementation.** The execution plan resembles
33
34
35
36 881 a branched-tree. In the first round, each process is assigned an approximately equal number of
37
38
39
40 882 files to merge locally. In the second round, even-numbered process retrieves the merged file of its
41
42
43
44 883 right adjacent process to merge with its local merged file. In the third round, processes whose ID
45
46
47
48 884 can be fully divided by four retrieve the merged file of its right adjacent process in the second
49
50
51
52 885 round and do the merging. This process continues recursively until all files are merged into a
53
54
55
56 886 single TPED file (round four).

57
58
59
60 887
61
62
63
64
65

1 888 **Figure 6. The scalability of Apache cluster-based schemas on input data size.** A. MapReduce
2
3
4
5 889 schema. B. HBase schema. C. Spark schema. As the number of input files increases from 10 to
6
7
8
9 890 186, the time costs of all three schemas with 12, 24 or 72 cores increase in a slower pace than that
10
11
12
13 891 of the input data size, especially when the number of cores is relatively large. The HBase schema
14
15
16
17 892 with 12 cores has the largest increase (from 375 to 5,479 seconds, ~14.6 fold).

18
19
20
21 893

22
23
24
25 894 **Figure 7. Comparing the strong scalability between traditional parallel/distributed methods**
26
27
28
29 895 **and Apache cluster-based schemas.** We fix the number of files at 93 and increase the number of
30
31
32
33 896 nodes/cores. The baseline for the parallel multiway-merge is one single core, while for the others
34
35
36
37 897 is one single node (4 cores). All methods/schemas show a degraded efficiency as computing
38
39
40
41 898 resources increase 16 fold from the baseline. Specifically, the efficiency of MapReduce-, HBase-,
42
43
44
45 899 Spark-based schemas drops to 0.83, 0.63 and 0.61 respectively, while the efficiency of parallel
46
47
48
49 900 multiway-merge and HPC-based implementations drops to 0.06 and 0.53 respectively.

50
51
52 901

53
54
55
56 902 **Figure 8. Comparing the weak scalability between traditional parallel/distributed methods**
57
58
59
60 903 **and Apache cluster-based schemas.** We simultaneously increase the number of cores and input

1 904 data sizes while fixing the ratio of file/core (parallel multiway-merge) or file/node (all others) at
2
3
4
5 905 ten. The baseline is the same as in the test of strong scalability. All but the MapReduce-based
6
7
8
9 906 schema have degraded efficiency, among which the HPC-based implementation has the steepest
10
11
12
13 907 degradation. Specifically, when computing resource increases 16 fold from the baseline, the
14
15
16
17 908 efficiency of MapReduce-, HBase- and Spark-based schemas changes to 3.1, 0.87 and 0.75
18
19
20
21 909 respectively, and for parallel multiway-merge and HPC-based implementations, the efficiency
22
23
24
25 910 reduces to 0.42 and 0.35 respectively.
26
27
28
29 911
30
31
32
33 912 **Figure 9. The performance anatomy of cluster-based schemas on increasing input data size.**
34
35
36
37 913 The number of cores in these experiments is fixed at 48. Time costs of all phases of the three
38
39
40
41 914 schemas have a linear or sub-linear correlation with the input data size. **a)** MapReduce schema:
42
43
44
45 915 The two MapReduce phases have a comparable time cost, increasing 6.3- and 3.1-fold
46
47
48
49 916 respectively as the number of input files increases from 10 to 186. **b)** HBase schema: The time
50
51
52
53 917 spent in each phase increases 4.2-, 5.6- and 5.0-fold respectively as the number of input files
54
55
56
57 918 increases from 10 to 186. The bulk loading and exporting phases together take up more than 80%
58
59
60
61 919 of total time expense. **c)** Spark schema: The time cost increases 5.8-, 6.0- and 6.0-fold
62
63
64
65

1 920 respectively for the three stages as the number of input files increases from 10 to 186 files. Like
2
3
4
5 921 the HBase schema, the first two stages of the Spark schema together account for more than 80% of
6
7
8
9 922 the total time cost.

10
11
12
13 923

14
15
16
17 924 **Figure 10. Execution speed comparison among Apache cluster-based schemas and**

18
19
20
21 925 **traditional methods.** Firstly, we compare of the speeds of the three Apache schemas with that of

22
23
24
25 926 three traditional methods which are single-process multiway-merge, parallel multiway-merge and

26
27
28
29 927 HPC-based implementations. As the number of input files increases from 10 to 186, the speeds of

30
31
32
33 928 Apache cluster-based schemas improve much more significantly than traditional methods. The

34
35
36
37 929 numbers in the figures indicate the ratio of the time cost of each traditional method to that of the

38
39
40
41 930 fastest Apache cluster-based schema. Secondly, we compare the processing speed among the three

42
43
44
45 931 Apache cluster-based schemas which are comparable to each other regardless of the input data

46
47
48
49 932 size. The MapReduce schema performs the best in merging 10 and 186 files; The HBase schema

50
51
52
53 933 performs the best in merging 20, 40 and 60 files; The Spark schema performs the best in merging

54
55
56
57 934 93 files.

58
59
60
61 935

62
63
64
65

1 936 **Figure S1. Pseudocodes of the MapReduce schema.**

2
3
4
5 937

6
7
8
9 938 **Figure S2. Pseudocodes of the HBase schema.**

10
11
12
13 939

14
15
16
17 940 **Figure S3. Pseudocodes of the Spark schema.**

18
19
20
21 941

22
23
24
25 942 **Tables**

26
27
28
29 943 **Table 1. Performance comparisons between VCTools and Apache cluster-based schemas**

30
31
32
33
34
35
36
37
38

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

39
40 944

41
42
43
44 945 **Table 2. Pros and Cons of MapReduce, HBase and Spark schemas**

45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Schemas	Pros	Cons
MapReduce	<ul style="list-style-type: none">• Good for large input data size and sufficient computing resources.• Simple architecture and least overheads given sufficient computing	<ul style="list-style-type: none">• Merging is not incremental.• Much overheads when computing resources are limited

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

	<p>resources.</p> <ul style="list-style-type: none"> • Best parallelism • Good for one-time merging. • Performance is stable. 	
HBase	<ul style="list-style-type: none"> • Good for intermediate input data size (≥ 20 and ≤ 100 VCF files). • Supports incremental merging. • Supports On-Line Analytical Processing (OLAP). • Best storage efficiency. 	<ul style="list-style-type: none"> • Users must determine region number in advance. • Has most local I/O. • Complex performance tuning.
Spark	<ul style="list-style-type: none"> • Good for large input data size (> 100 VCF files) and relative limited computing resources. • Keeps intermediate results in memory and least local I/O. • Good for subsequent statistical analysis on merged results. 	<ul style="list-style-type: none"> • Possibly weakened data locality during loading. • Slight unstable performance when computing resources exceeds needs of input data size. • Actual execution plan is not transparent. • Complex performance tuning.

- 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

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

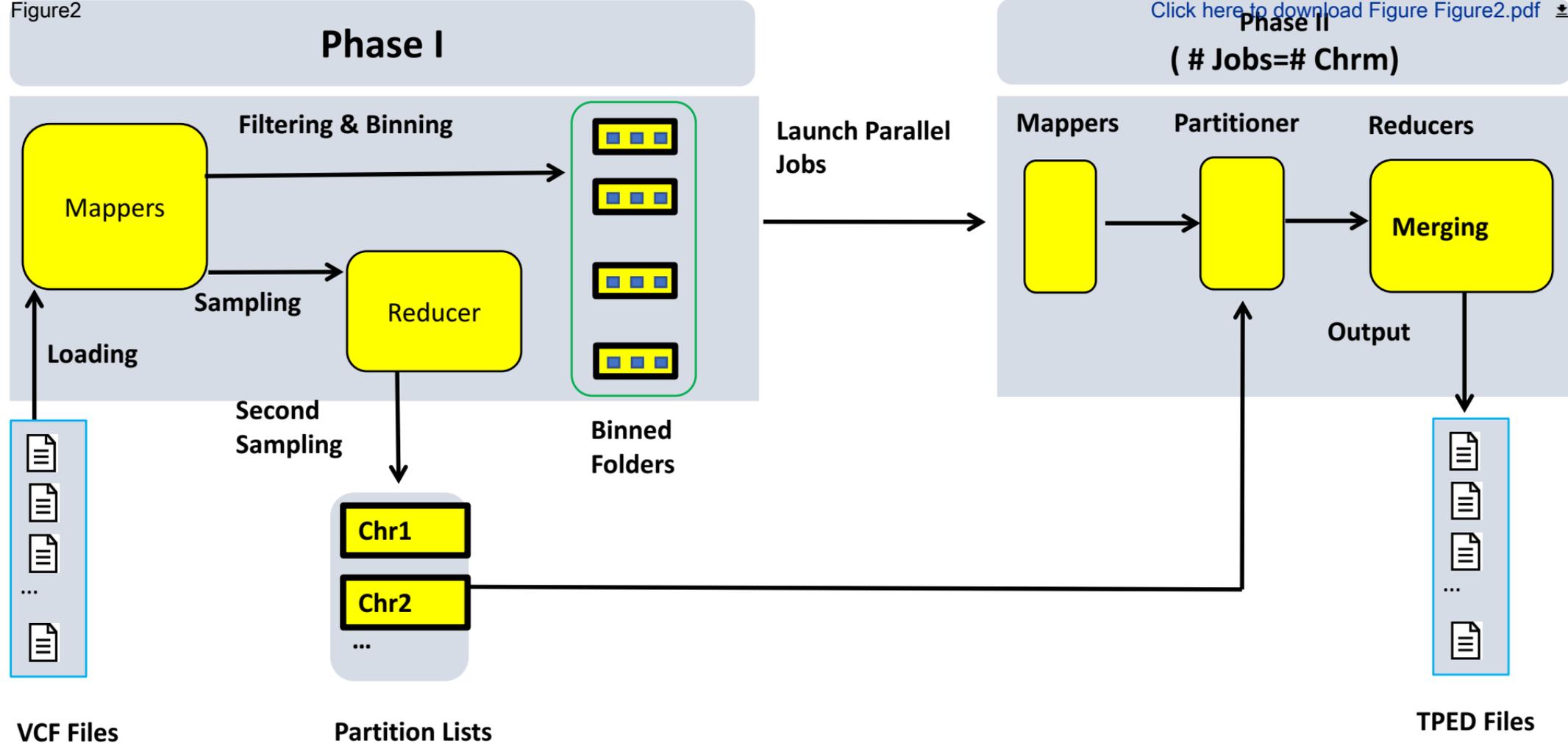


Genotypes

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](#)

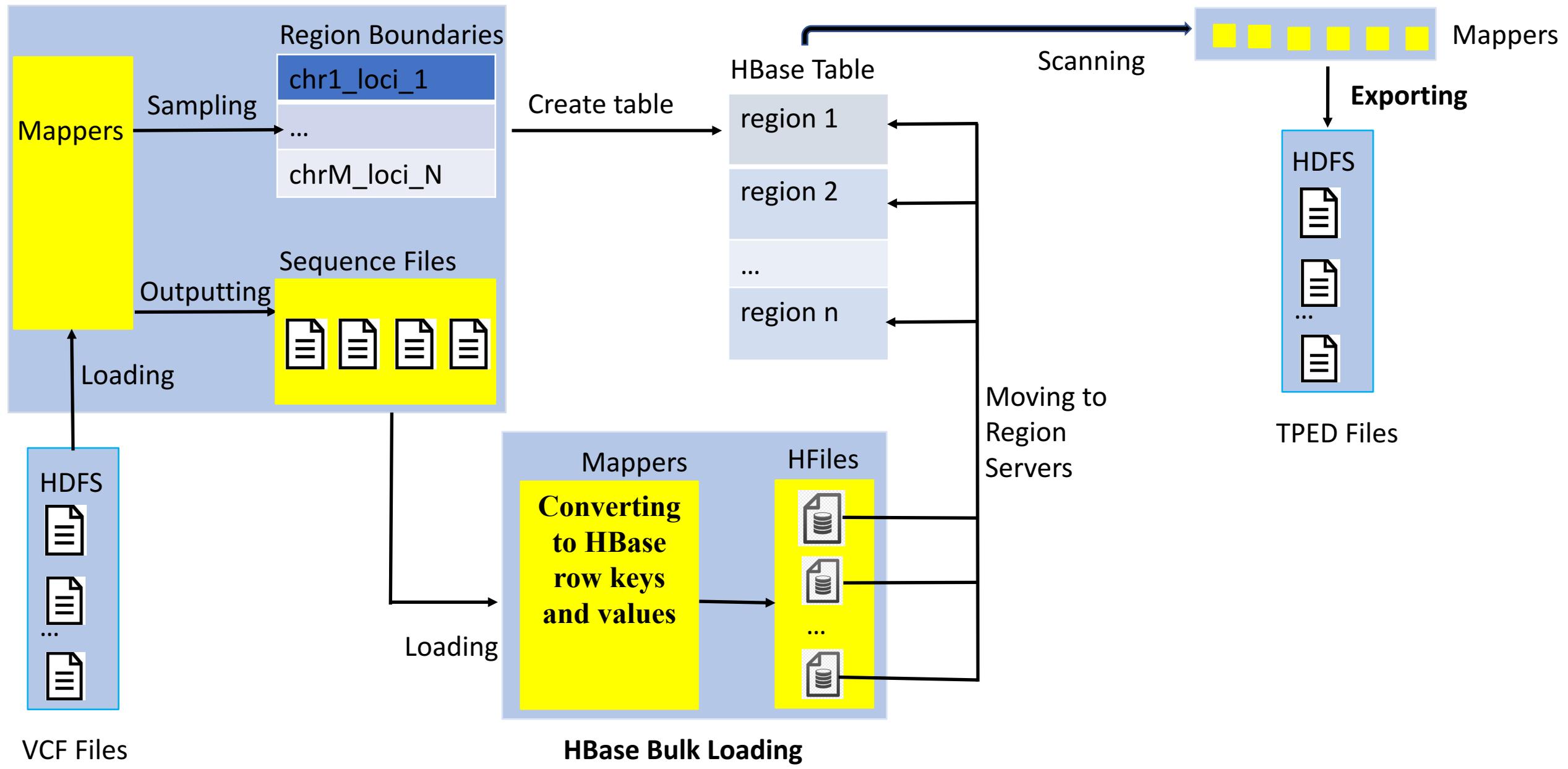
VCF Files

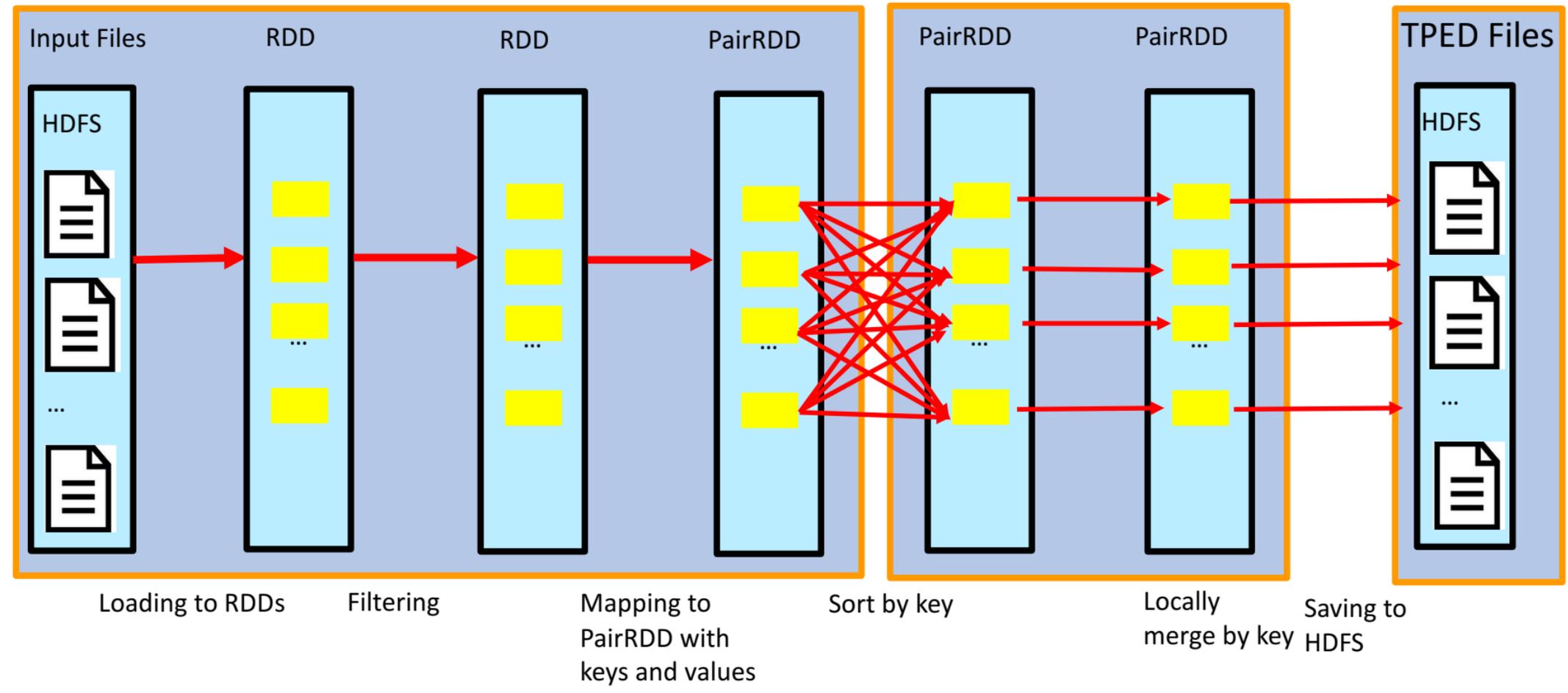
Partition Lists

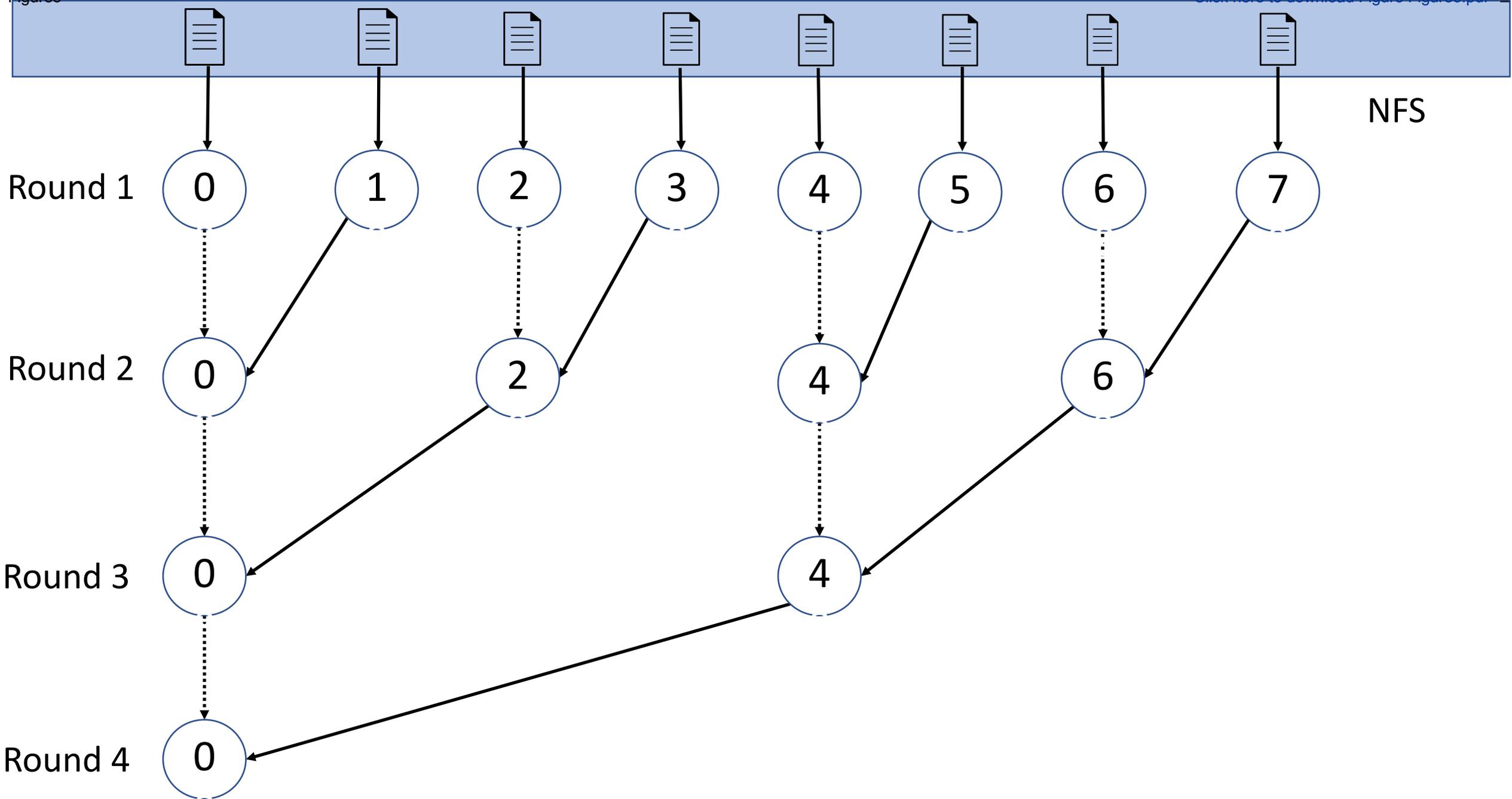
TPED Files

Figure3

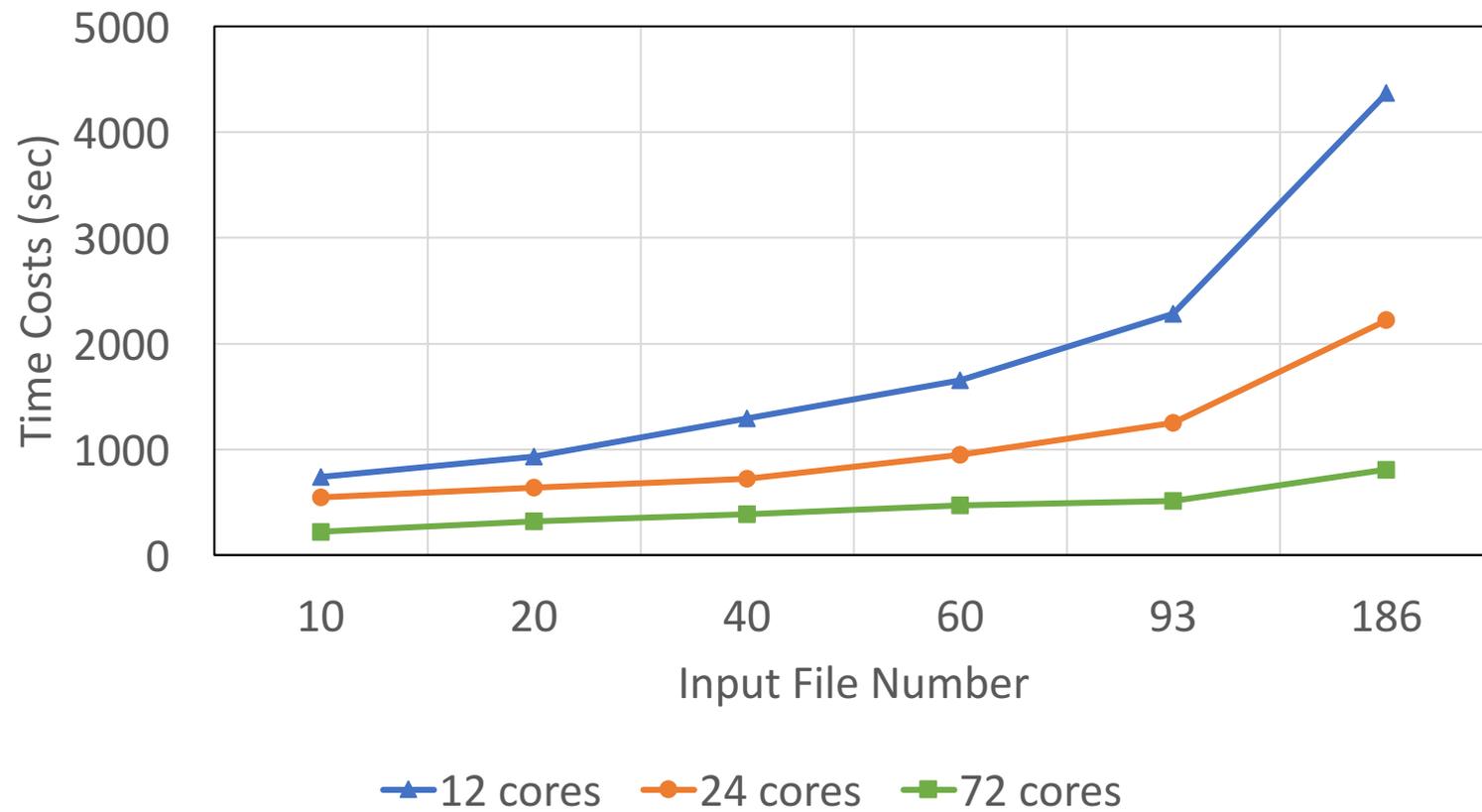
Sampling, Mapping & Filtering



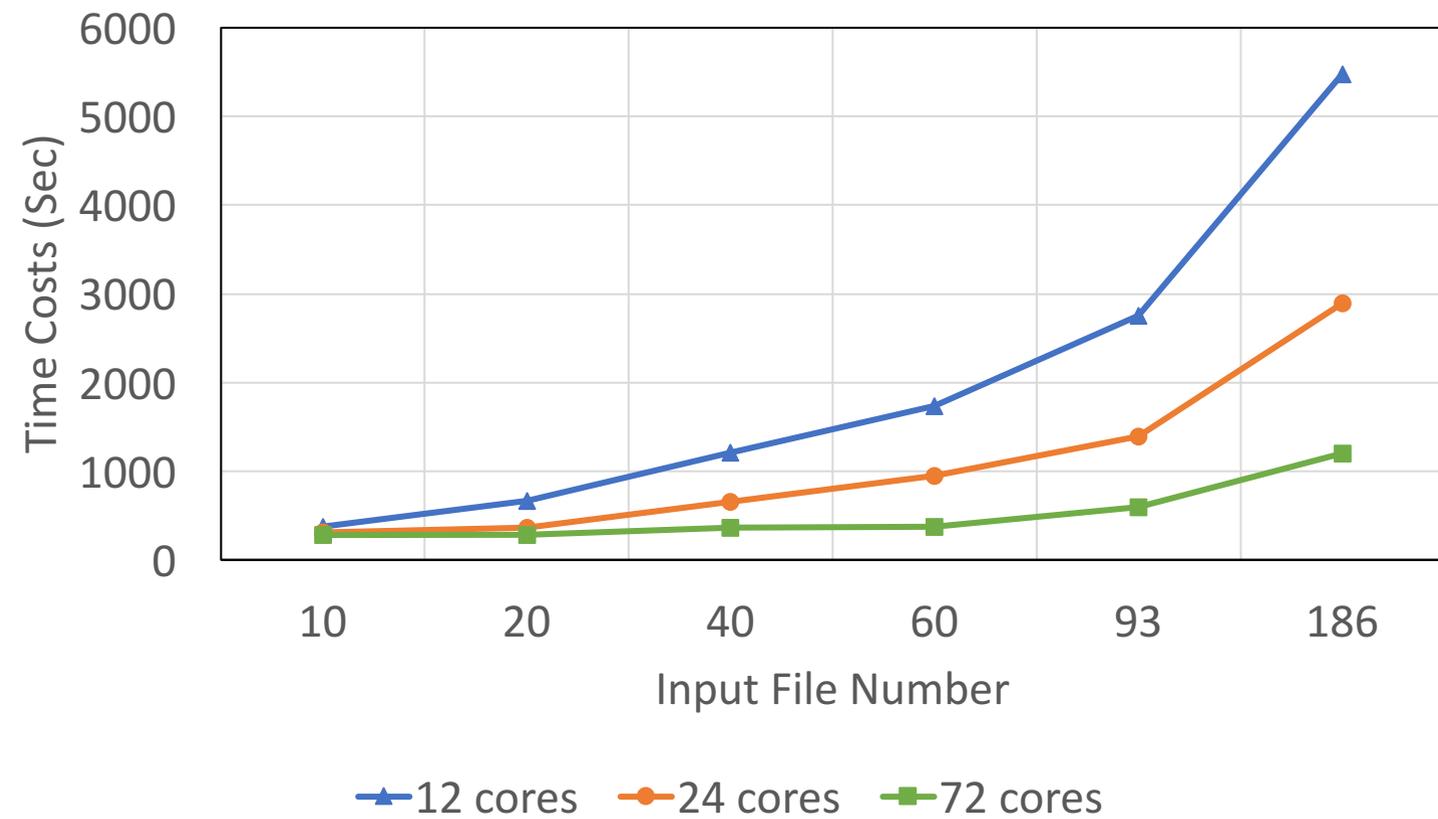




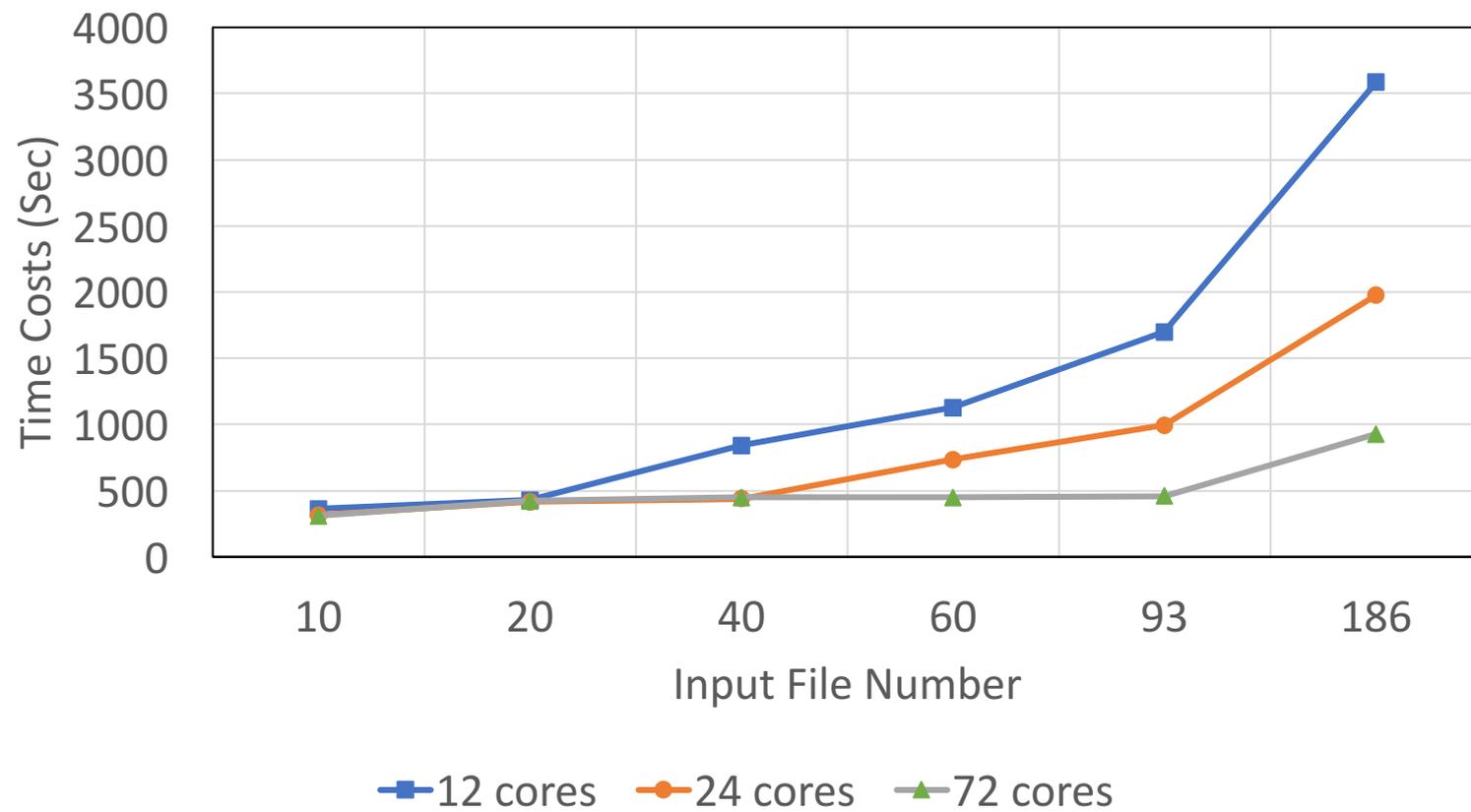
a) MapReduce Schema

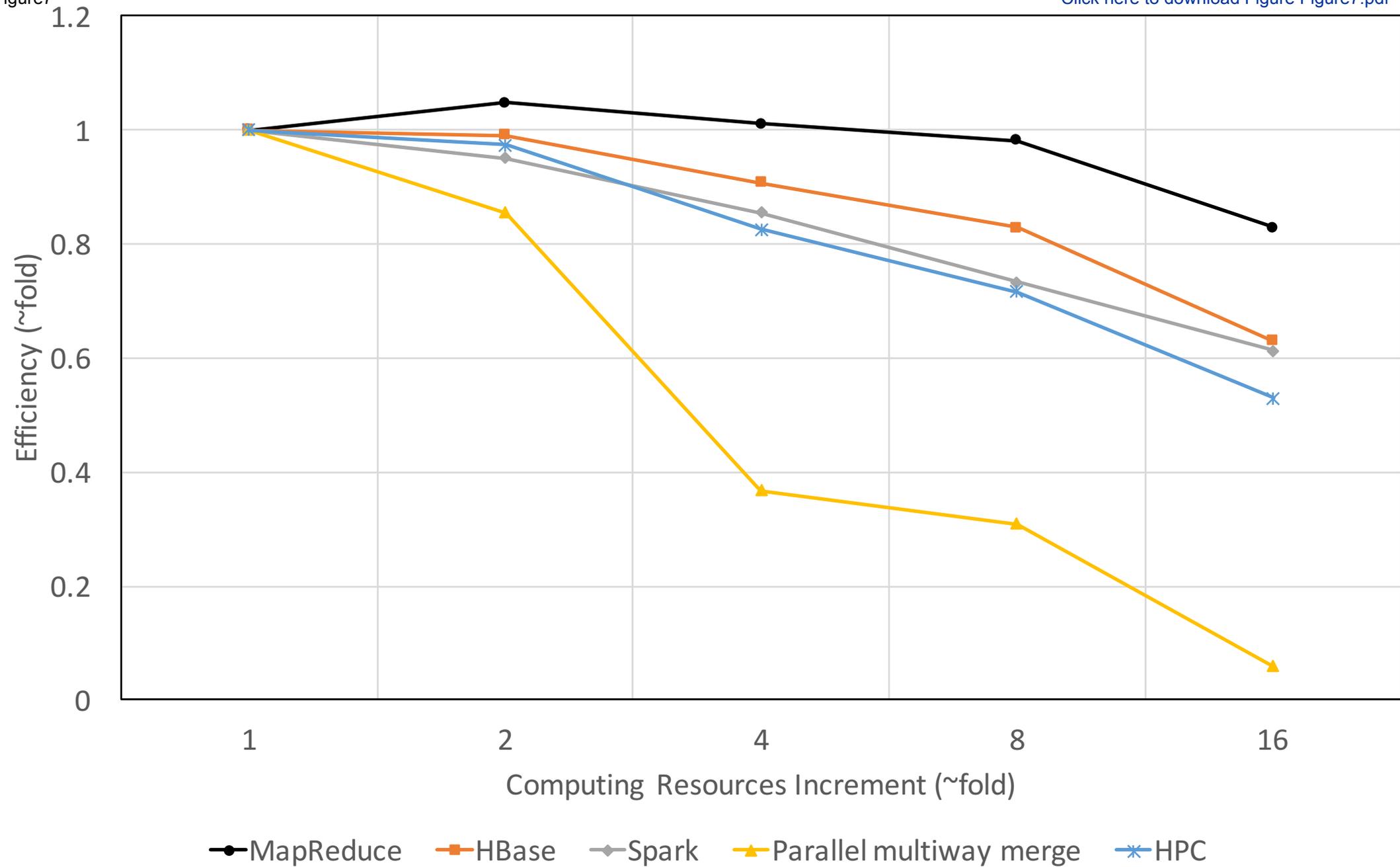


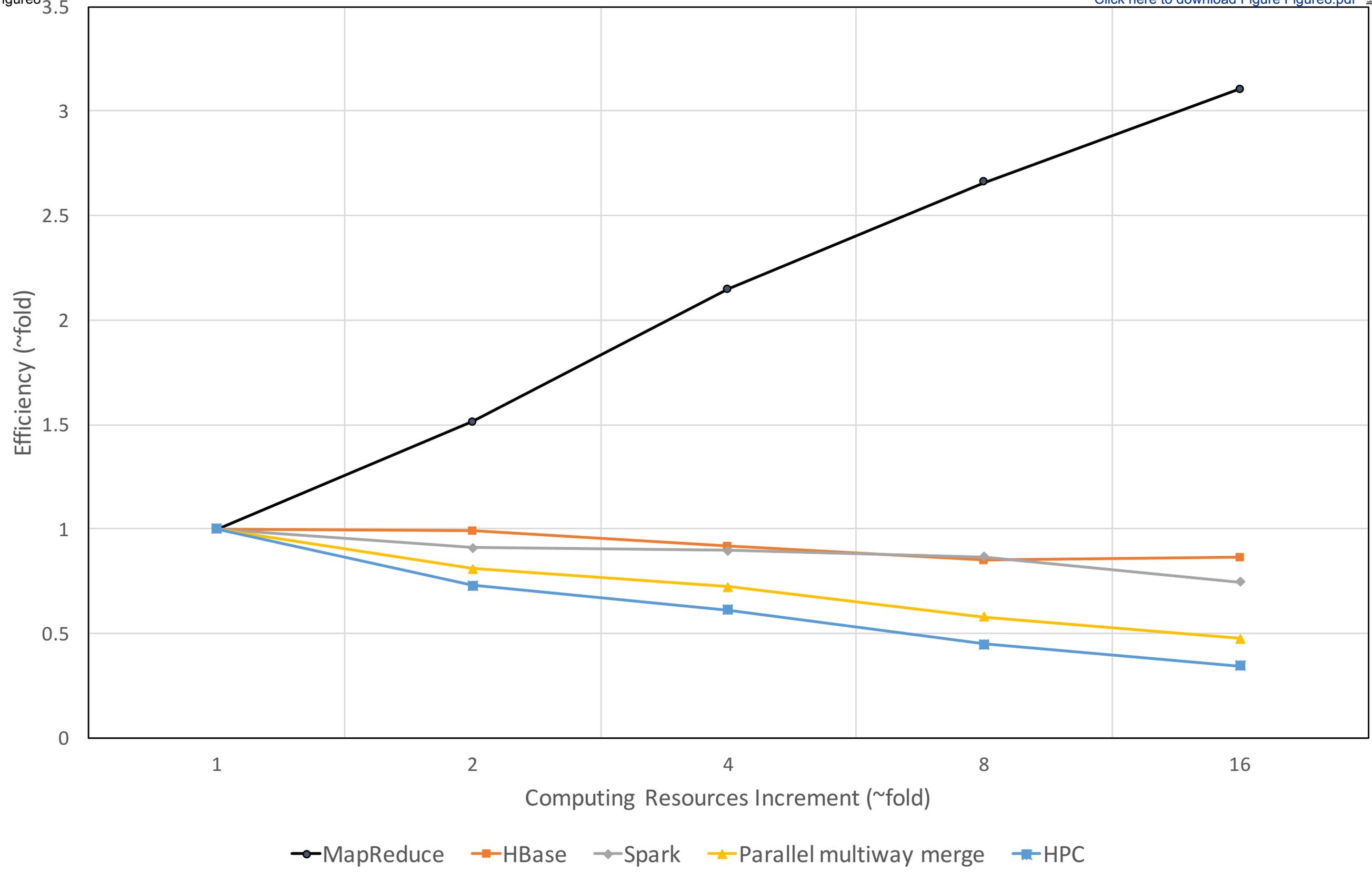
b) HBase Schema

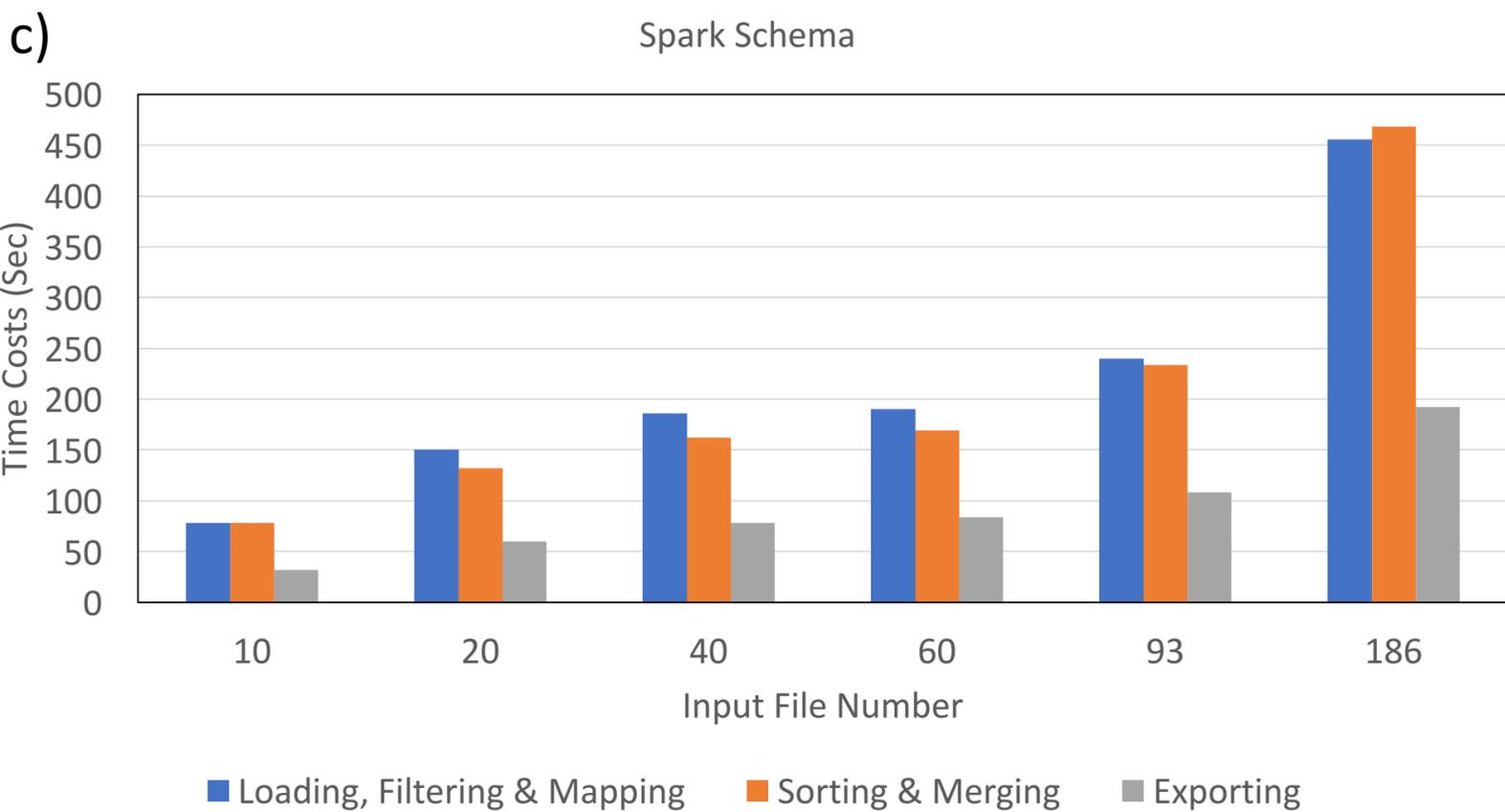
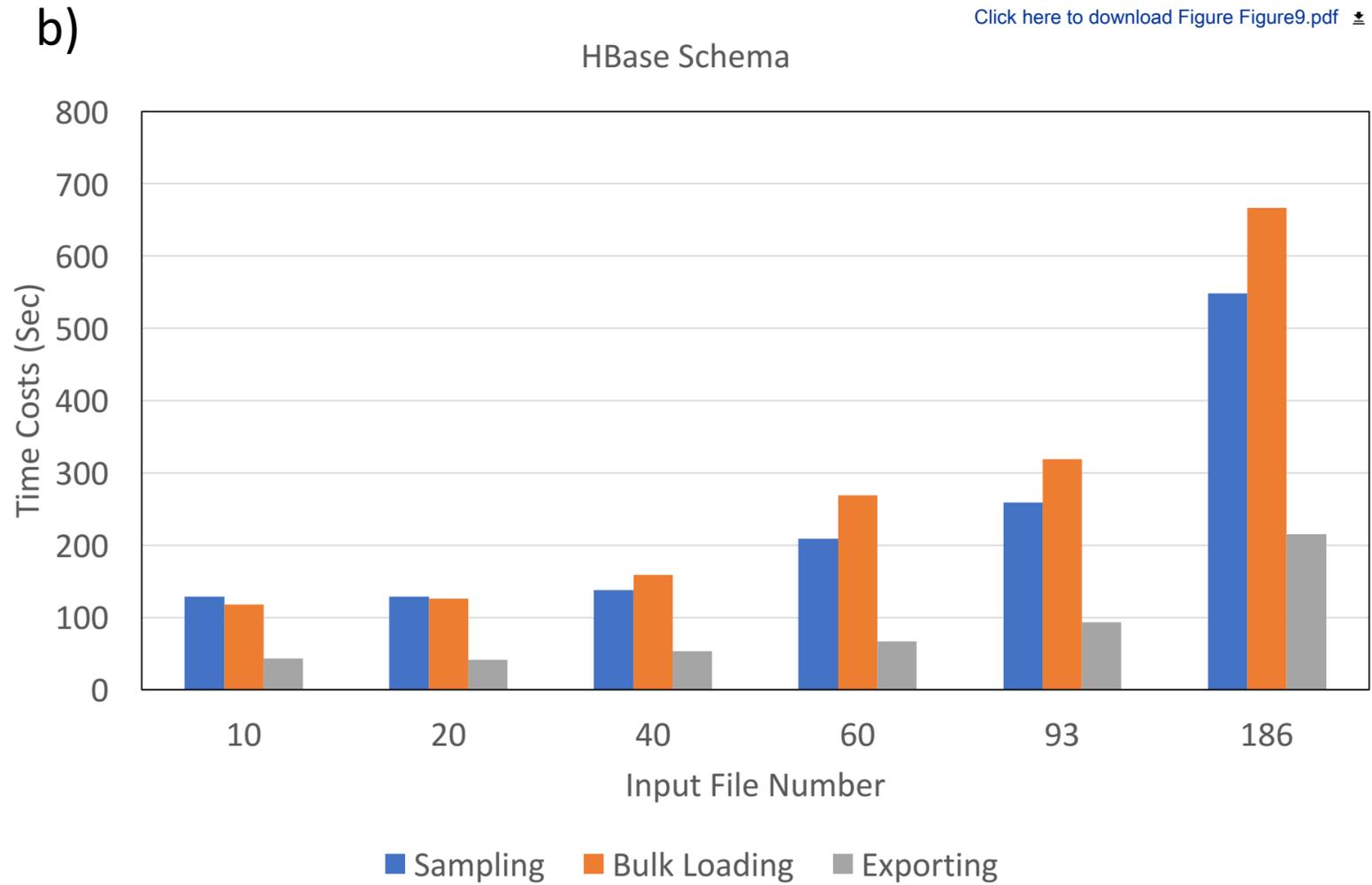
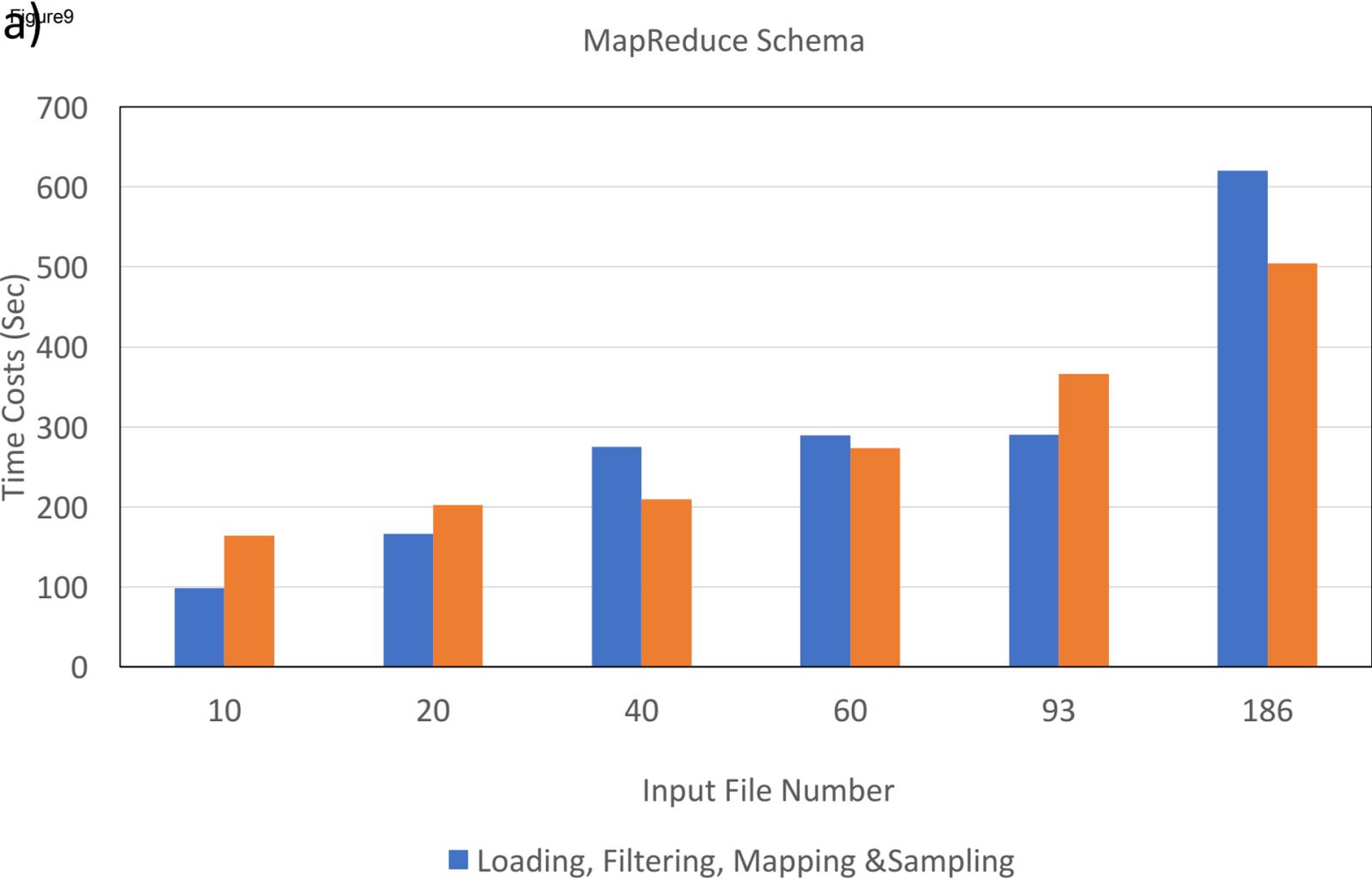


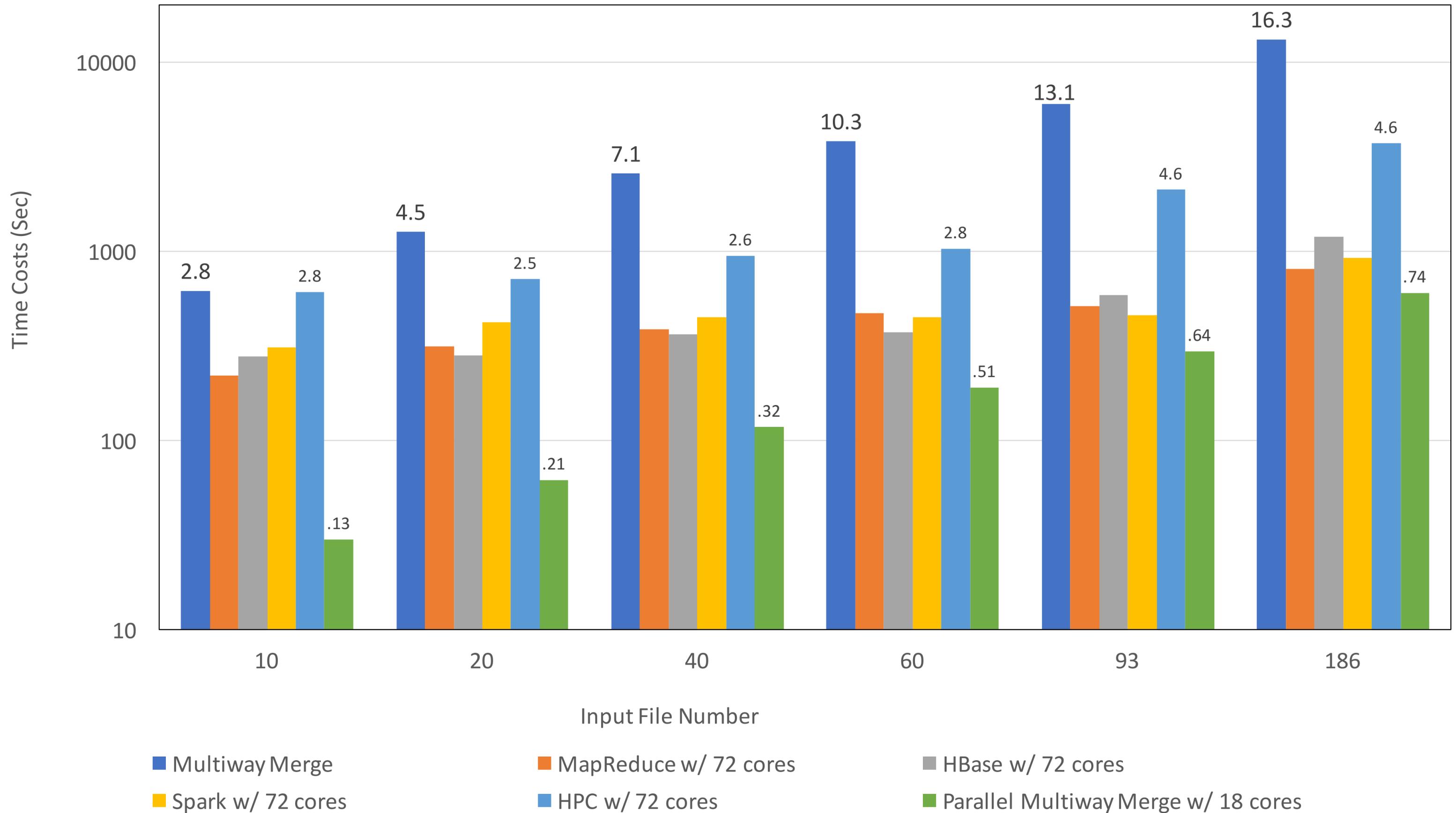
c) Spark Schema





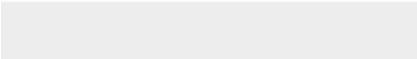
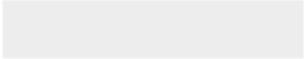








Click here to access/download
Supplementary Material
FigureS1.pdf





Click here to access/download
Supplementary Material
FigureS2.pdf



Click here to access/download
Supplementary Material
FigureS3.pdf



Click here to access/download
Supplementary Material
minor_revision_1.10_marked.docx

