

# GigaScience

## Bioinformatics Application on Apache Spark

--Manuscript Draft--

<b>Manuscript Number:</b>	GIGA-D-18-00131R1	
<b>Full Title:</b>	Bioinformatics Application on Apache Spark	
<b>Article Type:</b>	Review	
<b>Funding Information:</b>	National Key R&D Program of China (2017YFB0202600, 2016YFC1302500, 2016YFB0200400 and 2017YFB0202104)	Professor shaoliang peng
	National Natural Science Foundation of China (61772543, U1435222, 61625202, 61272056 and 61771331)	Professor shaoliang peng
	Guangdong Provincial Department of Science and Technology (2016B090918122)	Professor shaoliang peng
<b>Abstract:</b>	<p>With the rapid development of next-generation sequencing (NGS) technology, the ever-increasing genomic data pose a tremendous challenge to data processing. Therefore, there is an urgent need for highly scalable and powerful computational systems. Among the state-of-the-art parallel computing platforms, Apache Spark is a fast, general-purpose, in-memory, iterative computing framework for large-scale data processing, which ensures high fault tolerance and high scalability by introducing the resilient distributed dataset (RDD) abstraction. In terms of performance, Spark can be up to 100x faster in memory access and 10x faster in disk access than Hadoop. Moreover, it provides advanced APIs in Java, Scala, Python, and R. It also supports some advanced components, including Spark SQL for structured data processing, MLlib for machine learning, GraphX for graph computing, and Spark Streaming for stream computing. In this paper, we surveyed Spark-based applications in the NGS and other biological domains, such as epigenetics, phylogeny, and drug discovery. We believe that this survey provides a comprehensive guideline for bioinformatics researchers to apply Spark in their own fields.</p> <p>Keywords: next-generation sequencing; bioinformatics; Apache Spark; resilient distributed dataset; memory computing</p>	
<b>Corresponding Author:</b>	runxin guo  CHINA	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>		
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	runxin guo	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	runxin guo	
	yi zhao	
	xiangke liao	
	kenli li	
	quan zou	
	xiaodong fang	
	shaoliang peng	
<b>Order of Authors Secondary Information:</b>		

**Response to Reviewers:**

Dear Editor and Reviewers:

Thank you for your letter and for the reviewers' comments concerning our manuscript entitled "Bioinformatics Application on Apache Spark" (GIGA-D-18-00131). Those comments are all valuable and very helpful for revising and improving our paper, as well as the important guiding significance to our research. We have studied the comments carefully and have made correction which we hope meet with approval. Revised portion are marked in red in the paper. The main corrections in the paper and the responds to the editor and reviewers' comments are as following:

To Editor:

Comment 1: While one of the reviewers doesn't like the figures, another thinks you need more, and generally we side with the "more figures in a review the better". In particular, we would encourage improvements and more infographic like detail if you can. Improvements to the language and writing is also required, although if it passes review we will send it to a copy editor.

Response: Appreciate for your comment, taking into account the reviewers' comments and your opinions, we have improved the figures, removed unnecessary figures, and added a few necessary figures to help readers better understand the Spark framework and operating mechanism. In addition, we have improved my language and writing.

To Reviewer 1:

Comment 1: The paper would benefit from a \*slightly\* deeper description of the Spark architecture, in particular explaining the nature of DAGs and the way in which they permit optimizations. Also, some mention of the two deploy modes (where the driver program can either be run on the client machine, or on a worker node).

Response: Appreciate for your comment, in "THE SPARK FRAMEWORK" section, we have made a more detailed description of the Spark architecture, explained the nature of DAGs and the way in which they permit optimizations, introduced the two deploy modes: cluster mode and client mode.

Comment 2: The paper would also benefit from a section that examines the potential downsides of using Spark, for example the potential complexity in creating and maintaining a Spark cluster, and the learning curve involved in learning a new API and perhaps even language (especially given the Functional Programming nature of the API).

Response: Appreciate for your comment, in "DISCUSSION" section, we have discussed the disadvantages of Spark, including the applications that Spark is not suitable for, the complexity of creating and maintaining a Spark cluster, the time cost of large-scale input data from local to remote servers in slow networks, the complex learning curve.

Comment 3: With regards to style, there are a number of places in the paper where the definite article is used where it shouldn't, and vice versa. In the interest of readability and not distracting the reader, these should be addressed. A similar point can be made with regard to the over-use of certain prepositions (e.g. "besides"), which are called out in detail in the next section.

Response: Appreciate for your comment, we have addressed the use of the definite article in paper, and replaced "besides" with "in addition", "moreover", and "furthermore".

Comment 4: p.1 line 28: "data" is treated as a plural in the rest of the paper, therefore "pose" rather than "poses".

Response: Appreciate for your comment, we have changed "poses" to "pose".

Comment 5: p.1 line 34: "by introducing resilient distributed dataset" should be "by introducing the resilient distributed dataset" (i.e. use of definite article)

Response: Appreciate for your comment, we have changed "by introducing resilient distributed dataset" to "by introducing the resilient distributed dataset".

Comment 6: p.1 line 40: In the end, we discussed the challenges...and the future work...". I haven't found this discussion in the paper.

Response: Appreciate for your comment, we have deleted this sentence from "ABSTRACT" section, but added a "DISCUSSION" section to discuss the advantages and disadvantages of Spark, some issues to be considered about cloud computing in the future and other bioinformatics fields that have not yet been involved.

Comment 7: p.2 line 4: "MapReduce preforms" should be "MapReduce performs".

Response: Appreciate for your comment, we have changed "preforms" to "performs".

Comment 8: p.2 line 21: "introducing resilient distributed dataset" should "introducing the resilient distributed dataset".

Response: Appreciate for your comment, we have changed "introducing resilient distributed dataset" to "introducing the RDD abstraction".

Comment 9: p.2 line 38: The documentation of Spark describes the driver program as "The process running the main () function of the application and creating the SparkContext". It does not "deploy the Spark operating environment". Perhaps the authors meant "deploy TO the Spark operating environment" but even here this would be incorrect, as the spark-submit script does this.  
Response: Appreciate for your comment, we have rewritten this part to provide the correct description in "THE SPARK FRAMEWORK" section.

Comment 10: p.2, line 43: As well as Scala, Spark provides APIs in Java, Python and more recently R. This flexibility is important to researchers when deciding whether to use Spark or not.  
Response: Appreciate for your comment, we have mentioned that Spark provides advanced APIs in Java, Scala, Python and R in "THE SPARK FRAMEWORK" section.

Comment 11: p.2, line 58: It is questionable that "the most important feature of RDD" is the fault tolerance. Certainly, it is "an important feature".  
Response: Appreciate for your comment, we have changed this part to "RDD achieves fault tolerance through a notion of lineage..." in "THE SPARK FRAMEWORK" section.

Comment 12: p.3, line 13: The referenced image appears to be an \_example\_ of a spark task flow chart, rather than \_the general\_ Spark task processing flow. For the reader's sake, the paper should either describe what this particular task is doing (including the fact that it is reading and writing to HDFS in this case). Otherwise the reader may form incorrect opinions or simply be confused. Alternatively, drop the figure entirely.  
Response: Appreciate for your comment, we have updated this figure to show an example of how Spark computes job stages in "THE SPARK FRAMEWORK" section.

Comment 13: p.3, line 17: "Besides" as preposition. This is a little colloquial and has an additional "in any case" meaning. To avoid distracting the reader, consider replacing "besides" as a preposition with alternatives like "In addition", "Moreover", "Furthermore". This can be applied to the rest of the paper, and I won't call any more out by line number.  
Response: Appreciate for your comment, we have replaced "besides" with "in addition", "moreover", and "furthermore" in paper.

Comment 14: p.4, line 4: "Burrow-Wheeler aligner" - either "The Burrow-Wheeler aligner" or "Burrow-Wheeler alignment" read better  
Response: Appreciate for your comment, we have changed "Burrow-Wheeler aligner" to "The Burrow-Wheeler aligner".

Comment 15: p.4, line 19: "Results showed" - "The results showed"  
Response: Appreciate for your comment, we have changed "Result showed" to "The results showed".

Comment 16: p.4, line 32: "achieved the average speedup of" - "achieved an average speedup of"  
Response: Appreciate for your comment, we have changed "achieved the average speedup of" to "achieved an average speedup of".

Comment 17: p.6, line 58: Drop "And" from the start of the sentence.  
Response: We have dropped "And" from the start of the sentence.

Comment 18: p.7, line 1: "Experiments results" - "Experimental results"  
Response: Appreciate for your comment, we have changed "Experiments results" to "Experimental results".

Comment 19: p.7, line 4: Perhaps it's worth pointing out that this is an example of the platform itself suggesting a new algorithm, rather than simply re-implementing an existing algorithm on the new platform. Similarly, for line 19 of this page.  
Response: Appreciate for your comment, we have pointed out that these two are examples of the Spark platform itself suggesting new algorithms in "SPARK IN ASSEMBLY" section.

Comment 20: p.7, line 23: Is SA-BR-MR running on Hadoop? (I ask because MR is a valid algorithm on Spark as well).  
Response: Appreciate for your comment, SA-BR-MR is running on Hadoop according to the reference paper.

Comment 21: p.8, line 17: "Results..." - "The results..."  
Response: Appreciate for your comment, we have changed "Results" to "The results".

Comment 22: p.8, line 41: "noises" - "noise".  
Response: Appreciate for your comment, we have changed "noises" to "noise".

Comment 23: p.9, lines 23-39: The epigenetics example just calls out the advantage of parallelization compared to sequential processing. Was there a parallelized attempt, perhaps using Hadoop, that the Yu N et al paper could demonstrate a superiority to?

Response: Appreciate for your comment, we have reviewed lots of related papers, but did not find some parallelized attempts.

Comment 24: p.10, line 8: "Saprk" - "Spark"

Response: Appreciate for your comment, we have changed "Saprk" to "Spark".

Comment 25: p.10, line 15: The term "checkpointing" is not explained even in the body of the referenced paper (Harnie D et al) and is probably best dropped.

Response: Appreciate for your comment, we have dropped the term "checkpointing" from the sentence.

Comment 26: p.11, line 43: Key Points section: I would respectfully disagree with the following statement: "We introduce the Apache Spark framework in detail, helping researchers to understand its architecture, programming model and processing mechanism." I think the authors do a good job of firstly, giving an \*overview\* of Spark (notwithstanding earlier points about getting into more detail), but I don't think this paper is a \*detailed\* description of Spark, its architecture or its programming model. Indeed, I don't think it \*needs\* to be - the survey of \*how Spark has successfully been used\* is probably of primary interest to most readers. But it's best to be clear about the scope of the paper in the Key Points so as to set readers' expectations correctly.

Response: Appreciate for your comment, we have updated this key point to point out that we outline the Apache Spark framework to researchers to understand its architecture, programming model and processing mechanism, and we have made a more detailed description of the Spark architecture in "THE SPARK FRAMEWORK" section.

Comment 27: p.11, line 48: Key Points section: Similarly, to above, I would edit the third Key Point to set readers' expectations correctly. The paper in its current form does not include a "discussion on the future of parallel computing in bioinformatics" (and in my opinion it does not need to).

Response: Appreciate for your comment, considering your opinions and that of other reviewers, we have added a "DISCUSSION" section to discuss Spark's strengths, weaknesses, and challenges faced in this field.

To Reviewer 2:

Comment 1: While I agree, Spark has a lot of advantages over other parallel and distributed computing frameworks such as MapReduce, I feel the current tone and content are too one-sided. In my own experience, Spark is mostly only useful for processing very large amount of data. For smaller data sets, the scalability gained by Spark may not be enough to justify the up-front time required for setting up and configuring a Spark-enabled system. Also, there is no discussion on computing hardware requirement (local computer cluster or commercial cloud computing platforms), and issues related to transfer of large data sets over the Internet. All these issues need to be discussed.

Response: Appreciate for your comment, we have discussed the strengths and weaknesses of Spark and issues related to transfer of large data sets over the Internet on local computer cluster or commercial cloud computing in "DISCUSSION" section, and pointed out the official proposal for hardware requirements in "THE SPARK FRAMEWORK" section.

Comment 2: The sections on 'Spark in motif analysis' and 'Spark in genomic inference' are poorly written. The terms 'motif' and 'genomic inference' are not properly defined. Do they mean transcription factor binding motifs, or simply frequently occurring DNA sequence some defined regions in the genome (e.g., promoters, enhancers, etc.)? Also, the term 'genomic data inference' is not well defined. Presumably the authors are referring to inference in a population genomics context.

Response: Appreciate for your comment, here, motif refers to transcription factor binding sites (TFBS), genomic inference refers to the inference in population genomics text. We have updated these two sections in the paper to provide correct descriptions about the terms 'motif' and 'genomic inference'.

Comment 3: The caption of all four figures are way too simple. In most cases, especially for complicated flow diagrams like Fig 1 and Fig 3, there is no explanation or description of the content of the figure. I believe the authors intends to illustrate the inner working of Spark using these figures. Nonetheless, the content of these figures is not explained in the caption nor the main text.

Response: Appreciate for your comment, considering your opinions and that of other reviewers, we have updated the figures, removed unnecessary figures, and added a few necessary figures to help readers better understand the Spark framework and operating mechanism. Moreover, we have added some explanations and descriptions of the content of the figures in the captions.

Comment 4: Despite the authors claiming they 'discuss the future of parallel computing in bioinformatics' (Key Points #3), the manuscript barely talks about the future, other than saying 'Spark will provide promising performance for biological researchers in the future' (Conclusions). I think this is a lost opportunity. What do the authors see as the major limitations in the field at the moment? New hardware? Better integration with cloud computing platforms? New application areas, such as proteomics, metabolomics, biomedical text, electronic medical health record, etc.?

Response: Appreciate for your comment, we have updated the Key Points #3 and discussed some issues to be considered about cloud computing in the future and pointed out other bioinformatics fields that have not yet been involved.

Comment 5: In multiple occasions, the authors use 'And' to begin a sentence. I do not think it is grammatically correct.

Response: Appreciate for your comment, we have dropped "And" from all related sentences in paper.

Comment 6: Throughout the manuscript, author names are often cited as ([last name] [first initials]), but sometimes they are cited as ([last name] [all initials]) or ([last name] [first name]). Please make sure names are formatted consistently.

Response: Appreciate for your comment, we have updated author names as ([last name] [first initial]) in paper to make sure names are formatted consistently.

To Reviewer 3:

Comment 1: The manuscript needs to be rewritten to provide more practical information to bioinformatics research users how Apache Spark based bioinformatics tools are actively used in specific bioinformatics research domains and what are the advantages of using the Apache Spark.

Response: Appreciate for your comment, we have rewritten the paper to provide more practical information to bioinformatics research users how Apache Spark based bioinformatics tools are actively used in specific bioinformatics research domains and what are the advantages of using the Apache Spark. In "THE SPARK FRAMEWORK" section, we have made a more detailed description of the Spark architecture and the main abstraction RDD, explained the nature of DAGs and the way in which they permit optimizations, provided the official proposal for hardware requirements to help researchers better understand and use Spark. Moreover, we have added the "DISCUSSION" section to discuss the strengths and weaknesses of Spark, the applications that Spark is suitable for, and some issues must be considered. Researchers can comprehensively consider how to use Spark through these contents combined with biological issues in their field.

Comment 2: Table 1 shows basic APIs of Apache Spark. The reviewer cannot get a point why the authors included this table. Delete the table or keep the table with presenting how the APIs could be used in the bioinformatics tools.

Response: Appreciate for your comment, we have removed the Table 1 from the paper.

Comment 3: In the first paragraph of the Introduction section, "However, the existing bioinformatics tools cannot effectively handle such a large amount of data. In order to solve the issues, MapReduce, a programming model for parallel computation of large datasets, has been proposed [1]." sentences should be updated. MapReduce has not been proposed for bioinformatics tool. MapReduce framework was proposed for general purpose big data analysis in the distributed manner.

Response: Appreciate for your comment, we have updated these sentences to point out that MapReduce was proposed for processing large-scale datasets in a distributed manner in information technology rather than for bioinformatics tools.

Comment 4: The manuscript did not mention about "DataFrames" API that is an extension of RDD. Most of new Spark features use this new DataFrames, and it should be addressed in the manuscript.

Response: Appreciate for your comment, we have mentioned the two extensions of RDD: DataFrame and Dataset in "The SPARK FRAMEWORK" section. Users can seamlessly switch between the three through simple API calls.

Comment 5: Table 2 only shows the application domain, program name, URL, references. The table should be updated to include more meaningful informatics of tools such as pros/cons or specific features or the tools.

Response: Appreciate for your comment, we have updated the Table 2 to provide more meaningful informatics, including name, function, features, pros/cons and reference of applications.

Comment 6: Figures 1 to 4 are not necessary in the manuscript. Instead of these figures, the authors should consider how to express the relationship of bioinformatics



tools and Apache Spark effectively using figures.  
 Response: Appreciate for your comment, considering the opinions of editor and other reviewers, we have improved the figures, removed unnecessary figures, and added a few necessary figures to help readers better understand the Spark framework and operating mechanism. Figure 1 is mainly used to introduce Spark's cluster architecture, explain its main components and functions, and how the Spark application runs on the cluster. Figure 2 is mainly used to show some examples of narrow and wide dependencies to help readers understand RDD's dependencies. Because these dependencies are important in the splitting job stages of Spark, so we add this figure. Figure 3 is mainly used to show an example of how Spark computes job stages to help readers understand the RDD operating mechanism and DAG scheduling. So, we add this figure. Moreover, we carefully considered your comments about how to express the relationship of bioinformatics tools and Apache Spark effectively using figures. We feel that the relationship between bioinformatics tools and Spark cannot be well expressed in figures. Because the varied items in bioinformatics can cover almost all sorts of computational optimization problems. The Spark framework represents a possible solution for tasks in parallel computing with some certain characteristics. So, we tried our best to provide readers with some practical advices on using Spark framework based on computational characteristics of problems in "Discussion" section. In "DISCUSSION" section, we have discussed the strengths and weaknesses of Spark, the applications that Spark is suitable for, and some issues must be considered to help researchers to consider how to use Spark combined with biological issues in their field. We tried our best to improve the manuscript and made some changes in the manuscript. These changes will not influence the content and framework of the paper. And here we did not list the changes but marked in revise paper. We appreciate for Editor/Reviewers' warm work earnestly, and hope that the correction will meet with approval. Once again, thank you very much for your comments and suggestions.  
 Yours  
 Sincerely  
 Runxin GUO, Yi ZHAO, Xiangke LIAO, Kenli LI, Quan ZOU, Xiaodong FANG, Shaoliang PENG

<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No
<b>Experimental design and statistics</b>	Yes
Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a> . Information essential to interpreting the data presented should be made available in the figure legends.	
Have you included all the information requested in your manuscript?	
<b>Resources</b>	Yes
A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource</a>	

<p><a href="#">Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	Yes

# Bioinformatics Application on Apache Spark

Runxin GUO<sup>1†</sup>, Yi ZHAO<sup>3†</sup>, Xiangke LIAO<sup>1</sup>, Kenli LI<sup>2</sup>, Quan ZOU<sup>4\*</sup>, Xiaodong FANG<sup>5\*</sup>,  
Shaoliang PENG<sup>1,2\*</sup>

<sup>1</sup>College of Computer, National University of Defense Technology, Changsha 410073, China

<sup>2</sup>College of Computer Science and Electronic Engineering & National Supercomputer Centre in Changsha, Hunan University, Changsha 410082, China

<sup>3</sup>Institute of Computing Technology, Chinese Academy of Sciences, Beijing 100190, China

<sup>4</sup>School of Computer Science and Technology, Tianjin University, Tianjin 300350, China

<sup>5</sup>BGI Genomics, BGI-Shenzhen, Shenzhen 518083, China

[pengshaoliang@nudt.edu.cn](mailto:pengshaoliang@nudt.edu.cn) ; [zouquan@nclab.net](mailto:zouquan@nclab.net) ; [fangxd@bgitechsolutions.com](mailto:fangxd@bgitechsolutions.com)

\*: corresponding author; †: equal contributors

## ABSTRACT

With the rapid development of next-generation sequencing (NGS) technology, the ever-increasing genomic data pose a tremendous challenge to data processing. Therefore, there is an urgent need for highly scalable and powerful computational systems. Among the state-of-the-art parallel computing platforms, Apache Spark is a fast, general-purpose, in-memory, iterative computing framework for large-scale data processing, which ensures high fault tolerance and high scalability by introducing the resilient distributed dataset (RDD) abstraction. In terms of performance, Spark can be up to 100x faster in memory access and 10x faster in disk access than Hadoop. Moreover, it provides advanced APIs in Java, Scala, Python, and R. It also supports some advanced components, including Spark SQL for structured data processing, MLlib for machine learning, GraphX for graph computing, and Spark Streaming for stream computing. In this paper, we surveyed Spark-based applications in the NGS and other biological domains, such as epigenetics, phylogeny, and drug discovery. We believe that this survey provides a comprehensive guideline for bioinformatics researchers to apply Spark in their own fields.

**Keywords:** next-generation sequencing; bioinformatics; Apache Spark; resilient distributed dataset; memory computing

## INTRODUCTION

NGS technology has generated huge amounts of biological sequence data. In order to use these data efficiently, we need to store and analyze the data accurately and efficiently. However, the existing bioinformatics tools cannot effectively handle such a large amount of data. Therefore, there is an urgent need for scalable and powerful distributed computing tools to solve this problem. In the field of information technology, MapReduce [1] is a distributed parallel programming model and



1 methodology for processing large-scale datasets. It splits large-scale datasets into many key-value  
2 pairs through both the map and reduce phases, significantly improving performance and showing  
3 good scalability. By combining the Hadoop Distributed File System (HDFS) and MapReduce,  
4 Apache Hadoop can enable distributed processing of large amounts data in a reliable, efficient, and  
5 scalable way, where HDFS is mainly used for distributed storage of massive datasets and  
6 MapReduce performs distributed computing on these datasets. As a result, Hadoop has been adopted  
7 by the bioinformatics community in several areas [2], such as alignment [3], mapping [4] and  
8 sequence analysis [5].

9  
10  
11  
12  
13  
14  
15 However, due to its disk-based I/O access pattern, intermediate calculation results are not cached.  
16  
17 Therefore, Hadoop is only suitable for batch data processing, and shows poor performance for  
18 iterative data processing. To resolve this problem, Apache Spark [6] has been proposed, which is a  
19 faster general-purpose computing framework designed specifically to handle huge amounts of data.  
20  
21 Unlike Hadoop's disk-based computing, Spark performs memory computing by introducing the  
22 RDD abstraction. Since it is possible to store intermediate results in memory, it is more efficient for  
23 iterative operations. In terms of performance, Spark can be up to 100x faster in memory access than  
24 Hadoop [6]. Even if we compare between them based on the performance of the disk, the gap is  
25 more than 10 times [7]. In terms of flexibility, Spark provides high-level APIs in Java, Scala, Python,  
26 and R, and interactive shell. In terms of generality, Spark provides structured data processing,  
27 machine learning, graph computing, and stream computing capabilities by supporting some  
28 advanced components.

## 29 30 31 32 33 34 35 36 37 **THE SPARK FRAMEWORK**

38  
39 Spark is an open source cluster computing environment designed for large-scale data processing,  
40 developed by UC Berkeley AMP lab. It provides advanced APIs in Java, Scala, Python and R, and  
41 an optimized engine that supports general execution graphs. It also supports some advanced  
42 components, including Spark SQL for structured data processing, MLlib for machine learning,  
43 GraphX for graph computing, and Spark Streaming for stream computing.

44  
45  
46  
47  
48  
49 As in Figure 1, Spark application runs as independent processes on the cluster and are coordinated  
50 by the SparkContext in the driver program. There are two types of deploy modes depending on  
51 where the driver program is running: cluster mode and client mode. In the former, driver program  
52 runs on a worker node. In the latter, driver program runs on the client machine. First, SparkContext  
53 requests the executors on the worker nodes in the cluster from the cluster manager (either Spark's  
54 own Standalone cluster manager, Apache Mesos, or Hadoop YARN). These executors are processes  
55 that can run tasks and store data in memory or on disk for application. Next, the SparkContext will  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 send tasks to the executors to perform. Finally, the executors return the results to the SparkContext  
2 after the tasks are executed. In Spark, an application generates multiple jobs. A job is split into  
3 several stages. Each stage is a task set containing several tasks, which performs some calculations  
4 and produces some intermediate results. Task is the smallest unit of work in Spark, completing a  
5 specific thing on an executor.  
6  
7

8  
9 As the main abstraction in Spark, RDD is a read-only collection of objects partitioned on different  
10 nodes in the cluster so that the data in RDD can be processed in parallel. The data in RDD are stored  
11 in memory by default, but Spark automatically writes RDD data to disk if memory resources are  
12 low. RDD achieves fault tolerance through a notion of lineage [6], that is, if an RDD partition on a  
13 node is lost because of a node failure, the RDD automatically recalculates the partition from its own  
14 data source. Moreover, Spark provides two types of operations on RDD: transformation and action.  
15 The former defines a new RDD, and the latter returns a result or writes RDD data to the storage  
16 system. Transformation employs lazy operation [8], which means that the operation of generating  
17 another RDD from one RDD transformation is not executed immediately, and the calculation  
18 process is not actually started until an action is performed. Furthermore, each transformation  
19 operation generates a new RDD, the newly generated RDD depends on the original RDD. According  
20 to the different types of transformation operations, RDD's dependencies can be divided into narrow  
21 dependency and wide dependency. The former refers to that each partition in the generated RDD  
22 only depends on the parent RDD fixed partition, and the latter refers to the fact that each partition  
23 of the generated RDD depends on all partitions of the parent RDD. Figure 2 shows examples of  
24 narrow and wide dependencies. In addition, Spark also provides two extensions of RDD: DataFrame  
25 and Dataset. Spark users can seamlessly switch between the three through simple API calls.  
26  
27

28  
29 Furthermore, Spark adopts a directed acyclic graph (DAG) [9] to optimize execution process by  
30 splitting the submitted jobs into several stages according to the wide dependency. For narrow  
31 dependency, it divides related transformation operations into the same stage because they can  
32 perform pipelining operations and thus reduces the processing time of submitted jobs. Figure 3  
33 shows an example of how Spark computes job stages. In addition, if the partitions on a node are lost  
34 because of a node failure, Spark can utilize the DAG to recalculate the lost partitions.  
35  
36

37  
38 As for the hardware requirements, the official proposal is to have 4 to 8 disks per node, configure  
39 at least 8GB memory and 8 to 16 CPU cores per machine, and use a 10 Gigabit or higher network.  
40  
41

## 42 **SPARK IN ALIGNMENT AND MAPPING**

43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 The rapid development of NGS technology has generated a large amount of sequence data (reads),  
2 which has a tremendous impact on sequence alignment and mapping process. Currently, the  
3 sequence alignment and mapping process still consume a lot of time.  
4

5 The Smith-Waterman (SW) algorithm [10], which produces the optimal local alignment between  
6 two strings of nucleic acid sequences or protein sequences, is widely used in bioinformatics.  
7 However, SW algorithm requires a high computational cost due to high computational complexity.  
8 To speed up the algorithm, in 2015, Zhao G *et al* implemented the SW algorithm on Spark for the  
9 first time, called as SparkSW [11]. It consisted of three phases: data preprocessing, SW as map tasks  
10 and top K records as reduce tasks. Experimental results [11] showed that SparkSW was load-  
11 balancing and scalable with computing resources increased.  
12

13 However, SparkSW merely supports SW algorithm without the mapping location and traceback of  
14 optimal alignment, as a result, SparkSW executes slowly. Therefore, in 2017, Xu B *et al* proposed  
15 DSA [12], which employed Single Instruction Multiple Data (SIMD) instruction to parallel the  
16 sequence alignment algorithm at each worker node. Experimental results [12] showed that DSA  
17 achieved up to 201x speedup over SparkSW and almost linear speedup with the increase of cluster  
18 nodes.  
19

20 Subsequently, Xu B *et al* proposed CloudSW [13], an efficient distributed SW algorithm which  
21 leveraged Spark and SIMD instructions to accelerate the algorithm and provided APIs service in  
22 the cloud. Experimental results [13] showed that CloudSW achieved up to 3.29x speedup over DSA  
23 and 621x speedup over SparkSW. CloudSW also showed excellent scalability and achieved up to  
24 529 giga cell updates per second (GCUPS) in protein database search with 50 nodes in Aliyun.  
25

26 **The** Burrows-Wheeler aligner (BWA) is composed of BWA-backtrack [14], BWA-SW [15] and  
27 BWA-MEM [16] for performing sequence alignment and mapping in bioinformatics. Before the  
28 advent of Spark-based BWA tool, there were several other BWA tools based on big data technology,  
29 including BigBWA [17], Halvade [18] and SEAL [19]. However, they were based on Hadoop  
30 showing limited scalability and complex implementation.  
31

32 As a result, in 2015, Al-Ars Z *et al* [20] implemented three different versions of BWA-MEM and  
33 compared their performance: a native cluster-based version, a Hadoop version and a Spark version.  
34 Three implementations were evaluated on the same IBM Power7 and Intel Xeon servers with the  
35 WordCount example. **The results** [20] showed that simultaneous multithreading improved the  
36 performance of three versions of BWA-MEM, and the Spark version with 80 threads increased  
37 performance by up to 87% than the native cluster version using 16 threads. Furthermore, the Hadoop  
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

1 version with 4 threads increased performance by 17% and the Spark version with more threads  
2 increased performance by 27%.  
3

4 After then, in 2016, [Abuín J et al](#) proposed SparkBWA [21] which is composed of three main phases:  
5 the RDDs creation phase, the map phase, and the reduce phase. Experimental results [21] showed  
6 that for the BWA-backtrack algorithm, SparkBWA achieved an average speedup of 1.9x and 1.4x  
7 compared with SEAL and pBWA respectively. For the BWA-MEM algorithm, SparkBWA was  
8 1.4x faster than BigBWA and Halvade tools on average.  
9

10 However, SparkBWA required the data availability in the HDFS. In general, the input files were  
11 given in gzip format, which required first uncompressing the file before uploading it to the HDFS.  
12 Subsequently, this also slowed down the execution of BWA itself, since data on the HDFS had to  
13 be reformatted as appropriate input to the BWA program tasks running on the cluster. Finally, the  
14 output files produced by those BWA tasks required significant time to combine separately at the  
15 end.  
16

17 Therefore, in 2017, [Mushtaq H et al](#) employed Spark to propose StreamBWA [22], where the input  
18 data was being streamed directly from a compressed file. This file could either be located on the  
19 master node or on a URL, which eliminated the cost of execution time of downloading the file and  
20 then uncompressing it. Moreover, since the master node could stream data to the data nodes, the  
21 overhead of uploading data to the HDFS could also be hidden. The master node could also start  
22 combining the output files of BWA tasks running on the data nodes, in parallel, once they were  
23 available, further reducing the overall time. Experimental results [22] showed that this streaming  
24 distributed approach was approximately 2x faster than the non-streaming approach. Furthermore,  
25 StreamBWA was 5x faster than SparkBWA.  
26

27 Multiple sequence alignment (MSA) refers to the sequence alignment of three or more biological  
28 sequences, such as protein or nucleic acid sequences. One of representative tools for performing  
29 MSA is PASTA [23]. PASTA is a derivative of SATé [24], which produces highly accurate  
30 alignments in shared memory computers. However, PASTA is limited to processing small and  
31 medium datasets, because the computing power of shared memory systems cannot meet the memory  
32 and time requirements of large-scale datasets.  
33

34 Therefore, in 2017, [Abuín J et al](#) proposed PASTASpark [25], which allowed executions on a  
35 distributed memory cluster taking advantage of Spark. It employed an in-memory RDD of key-  
36 value pairs to parallel the calculating MSAs phase. Experiments were conducted on two different  
37 clusters (CESGA and AWS). The results [25] showed that PASTASpark achieved up to 10x  
38 speedups compared with single-threaded PASTA and was able to process 200,000 sequences in 24  
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

1 hours using only AWS nodes. Therefore, PASTASpark ensured scalability and fault tolerance which  
2 greatly reduced the time to obtain MSA.  
3

4 The probabilistic pairwise model [26] is widely used in all consistency-based MSA tools, such as  
5 MAFFT [27], ProbCons [28] and T-Coffee(TC) [29]. However, the global distributed memory  
6 cannot meet the ever-increasing sequence datasets, which causes the need of specialized distributed  
7 databases, such as HBase or Cassandra. As a result, in 2017, Lladós J *et al* employed Spark to  
8 propose a new tool, PPCAS [30], which could parallel the probabilistic pairwise model for large-  
9 scale protein sequences and store it in a distributed platform. Experimental results [30] showed that  
10 it was better with single node and provided almost linear speedup with the increase in the number  
11 of nodes. In addition, it could compute more sequences using the same memory.  
12  
13

14 NCBI BLAST [31, 32] is widely used to implement algorithms for sequence comparison. Before  
15 the Spark-based BLAST was created, several other BLAST tools had been proposed including  
16 mpiBLAST [33], GPU-BLAST [34] and CloudBLAST [35]. However, with the increasing number  
17 of genomic data, these tools showed limited scalability and efficiency.  
18

19 As a result, in 2017, Castro M *et al* proposed SparkBLAST [36], which utilized cloud computing  
20 and Spark framework to parallel BLAST. In SparkBLAST, Spark's *pipe* operator and RDDs were  
21 utilized to call BLAST as an external library and perform scalable sequence alignment. It was  
22 compared with CloudBLAST on both Google and Microsoft Azure Clouds. Experimental results  
23 [36] showed that SparkBLAST outperformed CloudBLAST in terms of speedup, scalability and  
24 efficiency.  
25

26 Metagenomics is crucial for studying genetic material directly from environmental samples.  
27 Fragment recruitment is the process of aligning reads to reference genomes in metagenomics data  
28 analysis. In 2017, Zhou W *et al* proposed MetaSpark [37], which employed Spark to recruit  
29 metagenomics reads to reference genomes.  
30

31 MetaSpark utilized the RDD of Spark to cache datasets in memory and scaled well along dataset  
32 size increments. It consisted of five steps including constructing k-mer RefindexRDD, constructing  
33 k-mer ReadlistRDD, seeding, filtering, and banded alignment. It was evaluated on a ten-node cluster  
34 working under the Spark standalone module where each node contained an 8-core CPU and 16 GB  
35 RAM. It employed about one million 75bp Illumina reads dataset and two references (the 194  
36 human gut genomes and the bacterial genomes) that were respectively 0.616GB and 1.3GB in size.  
37 Experimental results [37] showed that MetaSpark recruited more reads than FR-HIT [38] with the  
38 same parameters and 1 million reads. MetaSpark recruited 501,856 reads when there were 0.616  
39 GB human gut genome references, while FR-HIT recruited 489,638 reads. MetaSpark increased  
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

1 recruited reads by 2.5%. When references changed to a 1.3 GB bacterial genome, MetaSpark  
2 recruited 463,862 reads, while FR-HIT recruited 444,671 reads. MetaSpark increased recruited  
3 reads by 4%. Moreover, the results also showed that MetaSpark offered good scalability. Under a  
4 0.616 GB reference, run time for 0.1 million reads was 51 min under 4 nodes, and decreased slightly  
5 to 23.5 min under 10 nodes. For the 1 million read datasets, MetaSpark would crash under 4 nodes  
6 due to limited memory. Under 6 nodes, it finished running after 312 min and would sharply decrease  
7 to 201 min under 10 nodes.

## 13 SPARK IN ASSEMBLY

15 Due to short lengths of the NGS reads (<500 bp), they need to be assembled prior to further analysis,  
16 which is another important phase in sequence analysis workflow. In general, there are two types of  
17 assembly: the reference assembly and *de novo* assembly. The assembly algorithm includes two  
18 categories: overlap-layout-consensus (OLC) algorithm and the de Bruijn graph algorithm. The  
19 former is generally employed to assemble longer reads, while the latter shows a good performance  
20 in assembling short reads.

26 Before Spark-based distributed memory *de novo* assemblers were created, although there were some  
27 MPI-based assemblers (such as Ray [39], AbySS [40] and SWAP-Assembler [41]), they showed  
28 limited scalability, accuracy, and computational efficiency. Therefore, in 2015, Abu-Doleh A *et al*  
29 proposed Spaler [42] taking advantage of Spark and GraphX API. It consisted of two main parts:  
30 (a) de Bruijn graph construction, and (b) Contigs generating. **It** was evaluated with other MPI-based  
31 tools in terms of quality, execution time, and scalability. **Experimental** results [42] showed that  
32 Spaler had better scalability and it could achieve comparable or better assemble quality.

39 To resolve the large memory requirement problem of most OLC *de novo* assemblers, in 2017, Paul  
40 A *et al* [43] employed string graph reduction algorithms taking advantage of Spark. The proposed  
41 Spark algorithms were evaluated with a very large sample dataset. **The results** showed that this  
42 dataset was assembled by the proposed Spark algorithms using 15 virtual machines in 0.5 hours  
43 compared to the 7.5 hours of OLC based Omega [44] assembler.

49 **In addition, some new assembly algorithms have also been proposed based on the Spark platform  
50 itself.**

52 In 2016, Pan X *et al* [45] put forward a new assembling algorithm based on Spark which employed  
53 the method of matching K-2 bit to simplify the de Bruijn graph. This algorithm was evaluated using  
54 6 groups of DNA in the NCBI. Experimental results [45] showed that this strategy not only solved  
55 the problem of low efficiency based on the MapReduce algorithm, but also optimized the algorithm  
56 itself. The combination of these two aspects were very suitable for the large-scale DNA sequence  
57  
58  
59  
60  
61  
62  
63  
64  
65



1 assembling. Moreover, the results also showed that the new sequence assembling algorithm based  
2 on Spark could ensure accuracy of assembling results.

3  
4 To address the problem of poor assembling precision and low efficiency, in 2017, Dong G *et al* [46]  
5 proposed SA-BR-Spark, a new sequence assembly algorithm based on Spark. The authors first  
6 designed a precise assembling algorithm under the strategy of finding the source of reads based on  
7 the MapReduce and Eulerian path algorithm (SA-BR-MR). SA-BR-MR calculated 54 sequences  
8 which were randomly selected from animals, plants and microorganisms with base lengths from  
9 hundreds to tens of thousands from NCBI. All matching rates of 54 sequences were 100%. For each  
10 species, the algorithm also summarized the range of K which made the matching rates to be 100%.  
11 In order to verify the range of K value of hepatitis C virus (HCV) and related variants, the randomly  
12 selected eight HCV variants were calculated. The results confirmed the correctness of K range of  
13 hepatitis C and related variants from NCBI. After that, SA-BR-Spark was put forward. Experimental  
14 results [46] showed that SA-BR-Spark provided a superior computational speed compared with SA-  
15 BR-MR.  
16

## 17 **SPARK IN SEQUENCE ANALYSIS**

### 18 **Spark in variant analysis**

19 The GATK (Genome Analysis Toolkit) DNA analysis pipeline is widely used in genomic data  
20 analysis. Before Spark-based GATK tools were created, while several other tools had been  
21 developed to address the issue of scalability in the pipeline (such as Halvade [18] and Churchill  
22 [47]), they showed limited scalability, accuracy and computational efficiency.  
23

24 Therefore, in 2015, Mushtaq H *et al* [48] utilized Spark to propose a cluster-based GATK pipeline.  
25 To reduce the execution time, this approach kept data active in the memory between the map and  
26 reduce phases. By using runtime statistics of the active workload, it achieved a dynamic load  
27 balancing algorithm that could better utilize system performance. Experimental results [48] showed  
28 that this method achieved a 4.5x speedup compared to the multi-threaded GATK pipeline on a single  
29 node. In addition, when executed on a 4-node cluster, this approach was 63% faster than Halvade.  
30

31 After that, in 2016, Deng L *et al* proposed HiGene [49], which employed Spark to enable multi-  
32 core and multi-node parallelization of the GATK pipeline. HiGene put forward a dynamic  
33 computing resource scheduler and an efficient data skew mitigation method to improve performance.  
34 Experiments were conducted with the NA12878 whole human genome dataset. **The results** [49]  
35 showed that HiGene reduced the total running time from days to nearly an hour. Furthermore,  
36 compared with Halvade, HiGene was also 2x faster. Meanwhile, Li X *et al* employed Spark to  
37 propose GATK-Spark [50] to parallel the GATK pipeline by taking full account of compute,  
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

1 workload and I/O characteristics. It was built on top of ADAM format [51]. Experimental results  
2 [50] showed that GATK-Spark shortened the total running time from 20 hours to 30 minutes on 256  
3 CPU cores which achieved more than 37 times speedup.  
4

5 The advent of Spark provides the possibility of interactive processing for NGS data. In 2014,  
6 **Wiewiórka M et al** proposed SparkSeq [52] to build and run genomic analysis pipelines in an  
7 interactive way by using Spark. Experimental results showed that SparkSeq achieved 8.4–9.15 times  
8 speedup than SeqPig. Moreover, it could accelerate data querying up to 110x and reduce memory  
9 consumption by 13x.  
10

### 11 **Spark in motif analysis**

12 Due to the nature of NGS technology, the generated data are usually accompanied by some **noise**  
13 or other types of errors which are known as uncertain data [53]. Among these uncertain data, **there**  
14 **are some frequently recurring motifs, also known as transcription factor binding sites (TFBS)** [54].  
15 Mining motifs from these uncertain data is an important problem but a computationally intensive  
16 task. Before Spark-based mining algorithm was created, while several mining algorithms had been  
17 developed (such as HPSPM [55], DGSP [56] and SPAMC [57]), they showed limited scalability.  
18 Therefore, Jiang F *et al* [58] utilized Spark to propose a scalable algorithm for mining sequence  
19 motifs. This algorithm took advantage of Spark's RDDs and DAG, and allowed users to specify the  
20 minimum and maximum length of motif. Experiments were conducted with human genome datasets  
21 and bacteria DNA sequence datasets. **The results** [58] showed this approach could take a short  
22 period of time to extract accurate motifs.  
23

## 24 **SPARK IN OTHER BIOLOGICAL APPLICATIONS**

### 25 **Spark in population genomic inference**

26 Efficient score statistic methods [59] are widely applied in high-throughput **population genomic**  
27 **data inference**. A typical method of estimating the sampling distribution of the statistics is to employ  
28 asymptotic approximation, but it is inappropriate for small or uncommon variants. Although  
29 resampling methods [60] are appropriate for **the inference in population genomics context**, they  
30 greatly increase the computational burden of analysis. In order to tackle the computational challenge  
31 for resampling based inference, Bahmani A *et al* proposed SparkScore [61], a **distributed population**  
32 **genomic inference** approach taking advantage of Spark. SparkScore leveraged the nature of  
33 asymptotic and resampling inference based on efficient score statistics. Experiments on synthetic  
34 datasets using Amazon Elastic MapReduce (EMR) [61] demonstrated **the efficiency and scalability**  
35 **of SparkScore**.  
36

### 37 **Spark in epigenetics**

1 CpG islands (CGI) are important epigenetic markers, which play an essential role in epigenetics  
2 [62]. However, it is very challenging to investigate the CpG islands and their structures. Before  
3 Spark-based applications were developed, while several methods had been proposed to determine  
4 the CPG island (such as bisulfite modification-based methods), they were time-consuming and too  
5 costly. Thus, Yu N *et al* [63] utilized Spark to propose a novel CpG box model and a Markov model  
6 to redefine and investigate the CpG island which could greatly accelerate the analytic process.  
7 Experiments were conducted with Human and mouse chromosome sequences, 24 chromosomes and  
8 21 chromosomes. **The results** [63] showed this cloud-assisted method displayed considerable  
9 accuracy and faster processing power (6-7 times faster with 10 cores) compared with sequential  
10 processing.

### 11 **Spark in phylogeny**

12 Phylogeny reconstruction plays an important role in molecular evolutionary studies but faces  
13 significant computational challenges. Before Spark-based tools were created, while several tools  
14 had been put forward for phylogeny reconstruction, they could not scale well with a significant  
15 increase in data sets. Therefore, in 2016, Xu X *et al* proposed CloudPhylo [64], a fast and scalable  
16 Phylogeny reconstruction tool making use of Spark. It evenly distributed the entire computational  
17 workload among the working nodes. Experiment was conducted with the 5220 bacteria whole  
18 genome DNA sequences. **The results** [64] showed that CloudPhylo took 24508 seconds with one  
19 worker node and it could scale well as worker nodes increased. Moreover, CloudPhylo performed  
20 better than several existing tools when using more worker nodes. In addition, CloudPhylo achieved  
21 higher speedup on a larger dataset of about 100GB generated by simulation.

### 22 **Spark in drug discovery**

23 It is crucial to identify candidate molecules that affect disease-related proteins in drug discovery.  
24 Although the Chemogenomics project tries to identify candidate molecules using machine learning  
25 predictor programs [65-67], these programs spend a significant time and cannot be easily extended  
26 to multiple nodes. To migrate existing programs to multi-node clusters without changing the original  
27 programs, Harnie D *et al* proposed S-CHEMO [68] using **Spark**. In S-CHEMO, the intermediate  
28 data would be consumed again immediately on nodes that generated the data, reducing time and  
29 network bandwidth consumption. Experiments [68] compared S-CHEMO with the original pipeline,  
30 which showed almost linear speedup up to 8 nodes. **Moreover, this implementation also allowed**  
31 **easier monitoring.**

### 32 **Spark in Single-cell RNA sequencing**

1 Single-cell RNA sequencing (scRNA-seq) is crucial for understanding biological processes.  
2 Compared with standard bulk RNA-seq experiments, scRNA-seq experiments typically generate a  
3 greater number of cell profiles. Although there are already several RNA-seq processing pipelines  
4 (such as Halvade, SparkSeq and SparkBWA), they cannot process such a large number of profiles.  
5 Therefore, Falco [69] was created to process large-scale transcriptomic data in parallel by using  
6 Hadoop and Spark. Experiments were conducted with two public scRNA-seq datasets. **The results**  
7 [69] showed compared with a highly optimized single-node analysis, Falco was at least 2.6 times  
8 faster. Moreover, as the number of computing nodes increased, running time decreased.  
9 Furthermore, it allowed users to employ the low-cost spot instances of AWS which reduced the cost  
10 of analysis by 65%.

### 11 **Spark in variant association and population genetics studies**

12 Effectively analyzing thousands of individuals and millions of variants is a computationally  
13 intensive problem. Traditional parallel strategies such as MPI/OpenMP show poor scalability.  
14 While Hadoop provides an efficient and scalable computing framework, it is heavily dependent on  
15 disk operations. Therefore, in 2015, **O'Brien A et al** proposed VariantSpark [70] to parallel  
16 population-scale tasks based on Spark and associated machine learning library, MLlib. Experiments  
17 were conducted on 3000 individuals with 80 million variants, which showed that VariantSpark was  
18 80% faster than ADAM, Hadoop/Mahout implementation and ADMIXTURE [71]. Moreover,  
19 compared with R and Python implementations it was more than 90 % faster. And in 2017, **Di Z et**  
20 **al** proposed SEQSpark [72] to perform rare variant association analysis by using Spark. It was  
21 evaluated with whole-genome and simulated exome sequence data. The former was completed in  
22 1.5 hours and the latter in 1.75 hours. Moreover, it was always faster than Variant Association Tools  
23 and PLINK/SEQ, and in some cases running time was reduced to one percent.

### 24 **Spark in other works**

25 Biological simulations and experiments produce a large number of numerical datasets, and in 2017  
26 Klein M *et al* proposed Biospark [73] to process these data. Biospark was based on Hadoop and  
27 Spark, consisting of a set of Java, C++ and Python libraries. In addition, it provided the abstractions  
28 for parallel analysis of standard data types, including multidimensional arrays and images. To help  
29 parallel analysis of some common datasets, it also provided APIs and file conversion tools,  
30 including Monte Carlo, molecular dynamics simulations and time-lapse microscopy.

## 31 **DISCUSSION**

32 **Spark is an in-memory iterative computing framework designed for large-scale data processing. It**  
33 **is suitable for applications that require iterative operations on specific datasets. The greater the**

1 amount of data, the higher the computational intensity and the greater the benefit. When the data  
2 volume is small but the computational intensity is large, the benefit is relatively small. In addition,  
3 Spark is also suitable for applications where the amount of data is not particularly large but real-  
4 time statistical analyses are required.  
5

6  
7 However, due to the nature of RDD, Spark is not suitable for applications that require asynchronous  
8 fine-grained update in execution, such as web service storage or incremental web crawlers and  
9 indexes. In addition, we need to consider the potential complexity of creating and maintaining a  
10 Spark cluster. Moreover, when Spark runs on a commercial cloud computing platform such as AWS,  
11 there is a certain delay in the transmission of large-scale datasets over the Internet. This issue does  
12 not exist when Spark runs on a local computer cluster. Furthermore, we need to learn a new API  
13 and perhaps even language (especially given the functional programming nature of the API).  
14

15 Although Spark has been applied in some areas of bioinformatics and has achieved good results,  
16 other areas have not yet been involved, such as proteomics, biomedical text, and metabolomics.  
17 Moreover, as cloud computing and some web servers become more and more available, some issues  
18 must be considered, such as the time cost of large amounts of input data from local to remote servers  
19 in slow networks, cloud computing fees, data security and privacy.  
20

21 Table 1 summarizes the bioinformatics tools and algorithms based on Apache Spark.  
22

## 23 CONCLUSION

24 With the rapid development of NGS technology, a large number of genomic data have been  
25 generated, which poses a great challenge to traditional bioinformatics tools. For this reason, we have  
26 summarized the relevant works about Spark in bioinformatics and made a guideline on this topic.  
27 First, we make a comparison between Spark and Hadoop, and then outline the Spark cluster  
28 architecture, programming model, and processing mechanism. After that, we survey Spark-based  
29 applications in the NGS and other biological domains. A researcher who wants to get involved in  
30 this field can have a general understanding of Spark in bioinformatics through our survey.  
31

32 In summary, Spark is a fast and general-purpose computing framework designed for large-scale  
33 data processing. It ensures high fault tolerance and high scalability by introducing the RDD  
34 abstraction and DAG scheduling. We believe that bioinformatics applications based on Spark will  
35 provide promising performance for biological researchers in the future.  
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

## Key Points

- Apache Spark not only gives researchers a possibility of achieving efficient, scalable and fault-tolerant computing performance, but also supports various system workloads such as batch processing, iterative, interactive and stream computing.
- We **outline** the Apache Spark framework to help researchers to understand its architecture, programming model and processing mechanism.
- We present Spark-based applications that can be employed in bioinformatics and **discuss the strengths and weaknesses of Spark and the challenges faced in this field.**

## COMPETING INTERESTS

The authors declare that they have no competing interests.

## FUNDING

This work was supported by National Key R&D Program of China [grant numbers 2017YFB0202600, 2016YFC1302500, 2016YFB0200400 and 2017YFB0202104]; National Natural Science Foundation of China [grant numbers 61772543, U1435222, 61625202, 61272056 and 61771331]; and Guangdong Provincial Department of Science and Technology [grant number 2016B090918122].

## REFERENCES

1. Dean J, Ghemawat S. **MapReduce: simplified data processing on large clusters.** *Communications of the ACM* 2008, **51**(1):107-113.
2. Zou Q, Li X-B, Jiang W-R *et al.* **Survey of MapReduce frame operation in bioinformatics.** *Briefings in bioinformatics* 2013, **15**(4):637-647.
3. Zou Q, Hu Q, Guo M *et al.* **HAlign: Fast multiple similar DNA/RNA sequence alignment based on the centre star strategy.** *Bioinformatics* 2015, **31**(15):2475-2481.
4. Nguyen T, Shi W, Ruden D. **CloudAligner: A fast and full-featured MapReduce based tool for sequence mapping.** *BMC research notes* 2011, **4**(1):171.
5. Nordberg H, Bhatia K, Wang K *et al.* **BioPig: a Hadoop-based analytic toolkit for large-scale sequence data.** *Bioinformatics* 2013, **29**(23):3014-3019.
6. Zaharia M, Chowdhury M, Franklin MJ *et al.* **Spark: Cluster computing with working sets.** *HotCloud* 2010, **10**(10-10):95.
7. Han Z, Zhang Y. **Spark: A Big Data Processing Platform Based on Memory Computing.** In: *Seventh International Symposium on Parallel Architectures, Algorithms and Programming: 2016.* 172-176.



- 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
8. Zaharia M, Chowdhury M, Das T *et al.* **Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing.** In: *Proceedings of the 9th USENIX conference on Networked Systems Design and Implementation: 2012.* USENIX Association: 2-2.
9. Convolbo MW, Chou J. **Cost-aware DAG scheduling algorithms for minimizing execution cost on cloud resources.** *Journal of Supercomputing* 2016, **72**(3):985-1012.
10. Smith TF, Waterman MS. **Identification of common molecular subsequences.** *Journal of Molecular Biology* 1981, **147**(1):195-197.
11. Zhao G, Ling C, Sun D. **SparkSW: Scalable Distributed Computing System for Large-Scale Biological Sequence Alignment.** In: *Ieee/acm International Symposium on Cluster, Cloud and Grid Computing: 2015.* 845-852.
12. Xu B, Li C, Zhuang H *et al.* **DSA: Scalable Distributed Sequence Alignment System Using SIMD Instructions.** In: *Ieee/acm International Symposium on Cluster, Cloud and Grid Computing: 2017.* 758-761.
13. Xu B, Li C, Zhuang H *et al.* **Efficient Distributed Smith-Waterman Algorithm Based on Apache Spark.** In: *IEEE International Conference on Cloud Computing: 2017.* 608-615.
14. Li H, Durbin R. **Fast and accurate short read alignment with Burrows–Wheeler transform:** Oxford University Press; 2009.
15. Li H, Durbin R. **Fast and accurate long-read alignment with Burrows–Wheeler transform.** *Bioinformatics* 2010, **26**(5):589-595.
16. Li H. **Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.** 2013, **1303.**
17. Abuín JM, Pichel JC, Pena TF *et al.* **BigBWA: approaching the Burrows–Wheeler aligner to Big Data technologies.** *Bioinformatics* 2015, **31**(24):4003.
18. Decap D, Reumers J, Herzeel C *et al.* **Halvade: scalable sequence analysis with MapReduce.** *Bioinformatics* 2015, **31**(15):2482-2488.
19. Pireddu L, Leo S, Zanetti G. **SEAL: a distributed short read mapping and duplicate removal tool.** *Bioinformatics* 2011, **27**(15):2159.
20. Al-Ars Z, Mushtaq H. **Scalability Potential of BWA DNA Mapping Algorithm on Apache Spark.** In: *SIMBig: 2015.* 85-88.
21. Abuín JM, Pichel JC, Pena TF *et al.* **SparkBWA: Speeding Up the Alignment of High-Throughput DNA Sequencing Data.** *Plos One* 2016, **11**(5):e0155461.
22. Alars HMA. **Streaming Distributed DNA Sequence Alignment Using Apache Spark.** 2017.

- 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
23. Mirarab S, Nguyen N, Warnow T. **PASTA: ultra-large multiple sequence alignment**. In: *International Conference on Research in Computational Molecular Biology: 2014*. Springer: 177-191.
  24. Liu K, Warnow TJ, Holder MT *et al*. **SATe-II: very fast and accurate simultaneous estimation of multiple sequence alignments and phylogenetic trees**. *Systematic biology* 2011, **61**(1):90-106.
  25. Abuín JM, Pena TF, Pichel JC. **PASTASpark: multiple sequence alignment meets Big Data**. *Bioinformatics* 2017, **33**(18):2948-2950.
  26. Miyazawa S. **A reliable sequence alignment method based on probabilities of residue correspondences**. *Protein Engineering* 1995, **8**(10):999.
  27. Katoh K, Standley DM. **MAFFT Multiple Sequence Alignment Software Version 7: Improvements in Performance and Usability**. *Molecular Biology & Evolution* 2013, **30**(4):772-780.
  28. Do CB, Mahabhashyam MS, Brudno M *et al*. **ProbCons: Probabilistic consistency-based multiple sequence alignment**. *Genome Research* 2005, **15**(2):330.
  29. Tommaso PD, Moretti S, Xenarios I *et al*. **T-Coffee: a web server for the multiple sequence alignment of protein and RNA sequences using structural information and homology extension**. *Nucleic Acids Research* 2011, **39**(Web Server issue):13-17.
  30. Lladós J, Guirado F, Cores F *et al*. **PPCAS: Implementation of a Probabilistic Pairwise Model for Consistency-Based Multiple Alignment in Apache Spark**; 2017.
  31. Altschul S, Gish W, Miller W *et al*. **Basic local alignment search tool**. *J. Mol. Biol.* 1990.
  32. C C, G C, V A *et al*: **BLAST+: architecture and applications**. *Bmc Bioinformatics* 2009, **10**(1):421.
  33. Darling AE, Carey L, Feng WC. **The design, implementation, and evaluation of mpiBLAST**. In.: Los Alamos National Laboratory; 2003.
  34. Vouzis PD, Sahinidis NV. **GPU-BLAST: using graphics processors to accelerate protein sequence alignment**. *Bioinformatics* 2010, **27**(2):182-188.
  35. Matsunaga A, Tsugawa M, Fortes J. **Cloudblast: Combining mapreduce and virtualization on distributed resources for bioinformatics applications**. In: *eScience, 2008 eScience'08 IEEE Fourth International Conference on: 2008*. IEEE: 222-229.
  36. Castro MRD, Tostes CDS, Dávila AMR *et al*. **SparkBLAST: scalable BLAST processing using in-memory operations**. *Bmc Bioinformatics* 2017, **18**(1):318.

- 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
37. Zhou W, Li R, Yuan S *et al.* **MetaSpark: a spark-based distributed processing tool to recruit metagenomic reads to reference genomes.** *Bioinformatics* 2017, **33**(7):1090-1092.
  38. Niu B, Zhu Z, Fu L *et al.* **FR-HIT, a very fast program to recruit metagenomic reads to homologous reference genomes.** *Bioinformatics* 2011, **27**(12):1704-1705.
  39. Boisvert S, Laviolette F, Corbeil J. **Ray: simultaneous assembly of reads from a mix of high-throughput sequencing technologies.** *Journal of Computational Biology A Journal of Computational Molecular Cell Biology* 2010, **17**(11):1519.
  40. Simpson JT, Wong K, Jackman SD *et al.* **ABYSS: a parallel assembler for short read sequence data.** *Genome Research* 2009, **19**(6):1117.
  41. Meng J, Wang B, Wei Y *et al.* **SWAP-Assembler: scalable and efficient genome assembly towards thousands of cores.** *Bmc Bioinformatics* 2014, **15**(S9):S2.
  42. Abu-Doleh A, Çatalyürek ÜV. **Spaler: Spark and GraphX based de novo genome assembler.** In: *IEEE International Conference on Big Data: 2015.* 1013-1018.
  43. Paul AJ, Lawrence D, Ahn TH. **Overlap Graph Reduction for Genome Assembly using Apache Spark.** In: *The ACM International Conference: 2017.* 613-613.
  44. Haider B, Ahn TH, Bushnell B *et al.* **Omega: an Overlap-graph de novo Assembler for Metagenomics.** *Bioinformatics* 2014, **30**(19):2717-2722.
  45. Pan X, Fu X-L, Dong G-F *et al.* **DNA sequence splicing algorithm based on Spark.** In: *Industrial Informatics-Computing Technology, Intelligent Technology, Industrial Information Integration (ICIICII), 2016 International Conference on: 2016.* IEEE: 52-56.
  46. Dong G, Fu X, Li H *et al.* **An Accurate Sequence Assembly Algorithm for Livestock, Plants and Microorganism Based on Spark.** *International Journal of Pattern Recognition & Artificial Intelligence* 2017, **31**(8).
  47. Kelly BJ, Fitch JR, Hu Y *et al.* **Churchill: an ultra-fast, deterministic, highly scalable and balanced parallelization strategy for the discovery of human genetic variation in clinical and population-scale genomics.** *Genome biology* 2015, **16**(1):6.
  48. Mushtaq H, Al-Ars Z. **Cluster-based Apache Spark implementation of the GATK DNA analysis pipeline.** In: *Bioinformatics and Biomedicine (BIBM), 2015 IEEE International Conference on: 2015.* IEEE: 1471-1477.
  49. Deng L, Huang G, Zhuang Y *et al.* **HiGene: A high-performance platform for genomic data analysis.** In: *IEEE International Conference on Bioinformatics and Biomedicine: 2016.* 576-583.

- 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
50. Li X, Tan G, Zhang C *et al.* **Accelerating large-scale genomic analysis with Spark.** In: *Bioinformatics and Biomedicine (BIBM), 2016 IEEE International Conference on: 2016.* IEEE: 747-751.
  51. Massie M, Nothaft F, Hartl C *et al.* **Adam: Genomics formats and processing patterns for cloud scale computing.** *EECS Department, University of California, Berkeley, Tech Rep UCB/EECS-2013-207* 2013.
  52. Wiewiórka MS, Messina A, Pacholewska A *et al.* **SparkSeq: fast, scalable and cloud-ready tool for the interactive genomic data analysis with nucleotide precision.** *Bioinformatics* 2014, **30**(18):2652-2653.
  53. Leung CK-S. **Uncertain frequent pattern mining.** In: *Frequent pattern mining.* Springer; 2014: 339-367.
  54. Das MK, Dai H-K. **A survey of DNA motif finding algorithms.** *BMC bioinformatics* 2007, **8**(7):S21.
  55. Shintani T, Kitsuregawa M. **Mining algorithms for sequential patterns in parallel: Hash based approach.** In: *Pacific-Asia Conference on Knowledge Discovery and Data Mining: 1998.* Springer: 283-294.
  56. Qiao S, Tang C, Dai S *et al.* **Partspan: Parallel sequence mining of trajectory patterns.** In: *Fuzzy Systems and Knowledge Discovery, 2008 FSKD'08 Fifth International Conference on: 2008.* IEEE: 363-367.
  57. Chen C-C, Tseng C-Y, Chen M-S. **Highly scalable sequential pattern mining based on mapreduce model on the cloud.** In: *Big Data (BigData Congress), 2013 IEEE International Congress on: 2013.* IEEE: 310-317.
  58. Jiang F, Leung CK, Sarumi OA *et al.* **Mining sequential patterns from uncertain big DNA in the spark framework.** In: *Bioinformatics and Biomedicine (BIBM), 2016 IEEE International Conference on: 2016.* IEEE: 874-881.
  59. Rao CR. **Large sample tests of statistical hypotheses concerning several parameters with applications to problems of estimation.** In: *Mathematical Proceedings of the Cambridge Philosophical Society: 1948.* Cambridge University Press: 50-57.
  60. Westfall PH, Young SS. **Resampling-based multiple testing: Examples and methods for p-value adjustment,** vol. 279: John Wiley & Sons; 1993.
  61. Bahmani A, Sibley AB, Parsian M *et al.* **SparkScore: leveraging apache spark for distributed genomic inference.** In: *Parallel and Distributed Processing Symposium Workshops, 2016 IEEE International: 2016.* IEEE: 435-442.

- 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
62. Erkek S, Hisano M, Liang C-Y *et al.* **Molecular determinants of nucleosome retention at CpG-rich sequences in mouse spermatozoa.** *Nature structural & molecular biology* 2013, **20**(7):868-875.
63. Yu N, Li B, Pan Y. **A cloud-assisted application over apache spark for investigating epigenetic markers on DNA genome sequences.** In: *Big Data and Cloud Computing (BDCloud), Social Computing and Networking (SocialCom), Sustainable Computing and Communications (SustainCom)(BDCloud-SocialCom-SustainCom), 2016 IEEE International Conferences on: 2016.* IEEE: 67-74.
64. Xu X, Ji Z, Zhang Z. **CloudPhylo: a fast and scalable tool for phylogeny reconstruction.** *Bioinformatics* 2016, **33**(3):438-440.
65. Wale N. **Machine learning in drug discovery and development.** *Drug Development Research* 2011, **72**(1):112-119.
66. Costello JC, Heiser LM, Georgii E *et al.* **A community effort to assess and improve drug sensitivity prediction algorithms.** *Nature biotechnology* 2014, **32**(12):1202-1212.
67. Sastry GM, Inakollu VS, Sherman W. **Boosting virtual screening enrichments with data fusion: coalescing hits from two-dimensional fingerprints, shape, and docking.** *Journal of chemical information and modeling* 2013, **53**(7):1531-1542.
68. Harnie D, Saey M, Vapirev AE *et al.* **Scaling machine learning for target prediction in drug discovery using apache spark.** *Future Generation Computer Systems* 2017, **67**:409-417.
69. Yang A, Troup M, Lin P *et al.* **Falco: a quick and flexible single-cell RNA-seq processing framework on the cloud.** *Bioinformatics* 2016, **33**(5):767-769.
70. O'Brien AR, Saunders NFW, Guo Y *et al.* **VariantSpark: population scale clustering of genotype information.** *Bmc Genomics* 2015, **16**(1):1-9.
71. Alexander DH, Novembre J, Lange K. **Fast model-based estimation of ancestry in unrelated individuals.** *Genome Research* 2009, **19**(9):1655.
72. Di Z, Zhao L, Li B *et al.* **SEQSpark: A Complete Analysis Tool for Large-Scale Rare Variant Association Studies Using Whole-Genome and Exome Sequence Data.** *American Journal of Human Genetics* 2017, **101**(1):115.
73. Klein M, Sharma R, Bohrer CH *et al.* **Biospark: scalable analysis of large numerical datasets from biological simulations and experiments using Hadoop and Spark.** *Bioinformatics* 2017, **33**(2):303-305.

**Table 1 Bioinformatics tools and algorithms based on Apache Spark**

Name	Function	Features	Pros/Cons	Reference
SparkSW	Alignment and mapping	Consists of three phases: data preprocessing, SW as map tasks and top K records as reduce tasks	Load-balancing, scalable, but without the mapping location and traceback of optimal alignment	[11]
DSA	Alignment and mapping	Leverages data parallel strategy based on SIMD instruction	Up to 201x speedup over SparkSW and almost linear speedup with the increase of cluster nodes	[12]
CloudSW	Alignment and mapping	Leverages SIMD instruction and provides APIs service in cloud	Up to 3.29x speedup over DSA and 621x speedup over SparkSW, high scalability and efficiency	[13]
SparkBWA	Alignment and mapping	Consists of three main stages: RDDs creation, map, and reduce phases, employs two independent software layers	For shorter reads, average 1.9x and 1.4x faster than SEAL and pBWA. For longer reads, average 1.4x faster than BigBWA and Halvade, but requires the data availability in HDFS	[21]
StreamBWA	Alignment and mapping	The input data are being streamed into the cluster directly from a compressed file	~2x faster than non-streaming approach, and 5x faster than SparkBWA	[22]
PASTASpark	Alignment and mapping	Employs an in-memory RDD of key-value pairs to parallel the calculating MSAs phase	Up to 10x speedup than single-threaded PASTA, ensures scalability and fault tolerance	[25]
PPCAS	Alignment and mapping	Based on the MapReduce processing paradigm in Spark	Better with single node and shows almost linear speedup with the increase of nodes	[30]
SparkBLAST	Alignment and mapping	Utilizes pipe operator and RDDs of Spark to call BLAST as an external library	Outperforms CloudBLAST in terms of speedup, scalability and efficiency	[36]
MetaSpark	Alignment and mapping	Consists of five steps: constructing k-mer RefindexRDD, constructing k-mer ReadlistRDD, seeding, filtering, and banded alignment	Recruits significantly more reads than SOAP2, BWA and LAST, and more reads by ~4 than FR-HIT, and shows good scalability and overall high performance	[37]
Spaler	Assembly	Employs GraphX API of Spark, consists of two main parts: de Bruijn graph construction and contigs generating	Shows better scalability and achieves comparable or better assemble quality than ABySS, Ray, and SWAP-Assembler	[42]
SA-BR-Spark	Assembly	Under the strategy of finding the source of reads based Spark platform	Shows a superior computational speed than SA-BR-MR	[46]
HiGene	Variant analysis	Puts forward a dynamic computing resource scheduler and an efficient data skew mitigation way	Reduces total running time from days to nearly an hour, and 2x faster than Halvade	[49]



GATK-Spark	Variant analysis	Takes full account of compute, workload, and I/O characteristics	Achieves more than 37 times speedup	[50]
SparkSeq	Variant analysis	Builds and runs genomic analysis pipelines in an interactive way by using Spark	Achieves 8.4-9.15 times speedup than SeqPig, and accelerate data querying up to 110x and reduce memory consumption by 13x	[52]
SparkScore	Population genomic inference	Employs asymptotic and resampling inference on the basis of efficient score statistics	Shows good efficiency and scalability with EMR on synthetic datasets with EMR	[61]
CloudPhylo	Phylogeny	Evenly distributes the entire workloads among the worker nodes	Shows good scalability and high efficiency, and Spark version is better than Hadoop version	[64]
S-CHEMO	Drug discovery	Intermediate data is immediately consumed again on the nodes that produce, saving time and bandwidth	Shows almost linear speedup up to 8 nodes compared with the original pipeline	[68]
Falco	Single-cell RNA sequencing	Consist of a splitting step, an optional pre-processing step and the main analysis step	At least 2.6x faster than a highly optimized single-node analysis, and with the increase of nodes, running time decreases	[69]
VariantSpark	Variant association and population genetics studies	Parallels population-scale tasks based on Spark and associated MLlib	80% faster than ADAM, Hadoop/Mahout version and ADMIXTURE, and more than 90% faster than R and Python implementations	[70]
SEQSpark	Variant association and population genetics studies	Splits large-scale datasets into many small blocks to perform rare variant association analysis	Always faster than Variant Association Tools and PLINK/SEQ, and in some cases, running time is reduced to one percent	[72]
BioSpark	data-parallel analysis on large numerical datasets	Consists of a set of Java, C++ and Python libraries, abstractions for parallel analysis of standard data types, some APIs and file conversion tools	Convenient, scalable, and useful, brings domain-specific features for biology	[73]

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

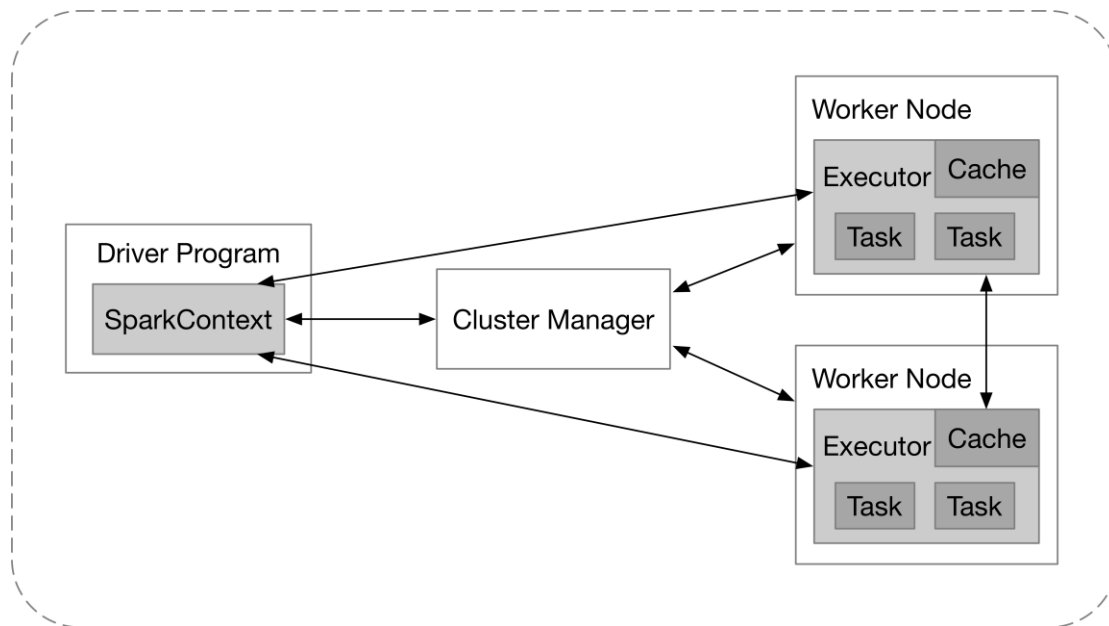


Figure 1: The cluster architecture of Spark

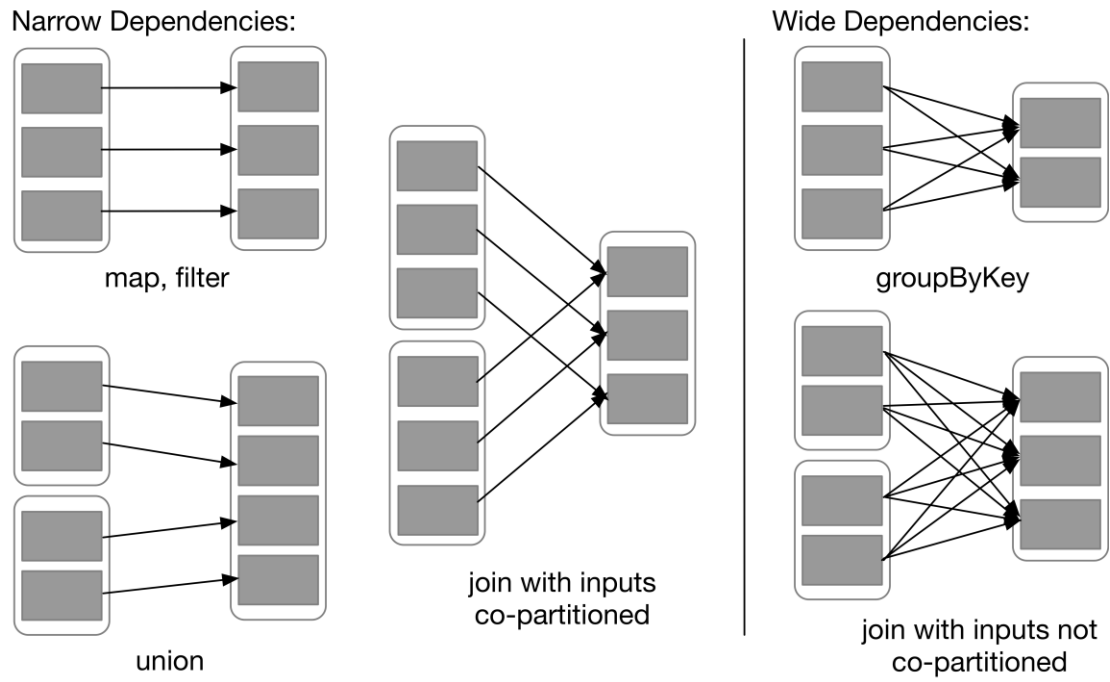


Figure 2: Examples of narrow and wide dependencies. Each box is an RDD, with partitions shown as shaded rectangles.

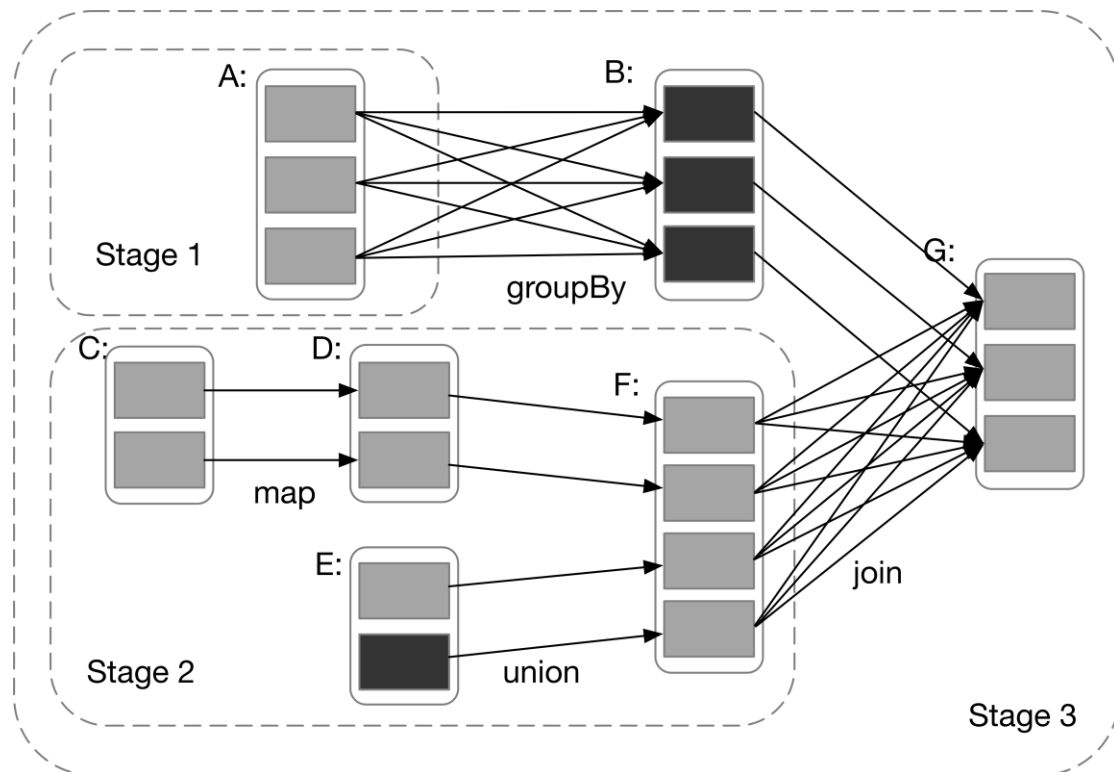


Figure 3: Example of how Spark computes job stages. Boxes with solid outlines are RDDs. Partitions are shaded rectangles, in black if they are already in memory. To run an action on RDD G, we build stages at wide dependencies and pipeline narrow transformation inside each stage. In this case, stage 1's output RDD is already in RAM, so we run stage 2 and then stage 3.



Click here to access/download  
**Supplementary Material**  
Cover Letter.docx





Click here to access/download  
**Supplementary Material**  
Home.md

