

BigDebug: Interactive Debugger for Big Data Analytics in Apache Spark

Muhammad Ali Gulzar, Matteo Interlandi, Tyson Condie, Miryung Kim
University of California, Los Angeles , USA
{gulzar, minterlandi, tcondie, miryung}@cs.ucla.edu

ABSTRACT

To process massive quantities of data, developers leverage data-intensive scalable computing (DISC) systems in the cloud, such as Google’s MapReduce, Apache Hadoop, and Apache Spark. In terms of debugging, DISC systems support post-mortem log analysis but do not provide interactive debugging features in realtime. This tool demonstration paper showcases a set of concrete use-cases on how BIGDEBUG can help debug Big Data Applications by providing interactive, realtime debug primitives. To emulate interactive step-wise debugging without reducing throughput, BIGDEBUG provides *simulated breakpoints* to enable a user to inspect a program without actually pausing the entire computation. To minimize unnecessary communication and data transfer, BIGDEBUG provides *on-demand watchpoints* that enable a user to retrieve intermediate data using a guard and transfer the selected data on demand. To support systematic and efficient trial-and-error debugging, BIGDEBUG also enables users to change program logic in response to an error at runtime and replay the execution from that step. BIGDEBUG is available for download at <http://web.cs.ucla.edu/~miryung/software.html>.

CCS Concepts

• **Software and its engineering** → **Software testing and debugging; Error handling and recovery; Development frameworks and environments;**

Keywords

Debugging, big data analytics, interactive tools, data-intensive scalable computing (DISC), fault localization and recovery

1. INTRODUCTION

An abundance of data in many disciplines of science, engineering, national security, health care, and business is now urging the need for developing Big Data Analytics. To process massive quantities of data in the cloud, developers leverage data-intensive scalable computing (DISC) systems such as Google’s MapReduce [4], Apache Hadoop [1], and Apache Spark [13]. In these DISC systems, scaling to large datasets is handled by partitioning data and

assigning tasks that execute a portion of the application logic on each partition in parallel. Unfortunately, this critical gain in scalability creates an enormous challenge for data scientists in understanding and resolving errors.

The application programming interfaces (API) provided by DISC systems expose a batch model of execution: applications are run in the cloud, and the results, including notification of runtime failures, are sent back to users upon completion. Therefore, debugging is done *post-mortem* and the primary source of debugging information is an execution log. However, the log presents only the *physical view*—the job status at individual nodes, the overall job progress rate, the messages passed between nodes, etc, but does not provide the *logical view*—which intermediate outputs are produced from which inputs, what inputs are causing incorrect results or delays, etc. Alternatively, a developer may test their program by downloading a small subset of big data from the cloud onto their local disk, and then run the application in local mode. However, this approach can easily miss errors, when the faulty data is not part of the downloaded subset.

We showcase BIGDEBUG, a library providing expressive and interactive debugging for big data analytics in Apache Spark [2]. This tool demonstration paper is based on our prior work on the design and implementation of interactive debugging primitives for Apache Spark [7]. Designing BIGDEBUG requires re-thinking the notion of breakpoints, watchpoints, and step-through debugging in a traditional debugger such as `gdb`. For example, simply pausing the entire computation would waste a large amount of computational resources and prevent correct tasks from completing, reducing the overall throughput. Requiring the user to inspect the millions of intermediate records at a watchpoint is also clearly infeasible.

BIGDEBUG provides **simulated breakpoints** that enable a user to inspect a program without actually pausing the entire computation. It also supports **on-demand watchpoints** that enable a user to retrieve intermediate data using a guard predicate and transfer the selected data on demand. To understand the flow of individual records within a data parallel pipeline, BIGDEBUG provides **data provenance** capability, which can help understand how errors propagate through data processing steps. To support systematic and efficient trial-and-error debugging, BIGDEBUG enables users to change program logic in response to an error at runtime through a **realtime code fix** feature and **selectively replay** the execution from that step. Under the maximum instrumentation settings, BIGDEBUG takes 2.5X more time than baseline Apache Spark. More results on performance overhead are described elsewhere [7].

This paper is organized as follows. Section 2 describes the background on Apache Spark. Section 3 describes individual debugging features of BIGDEBUG using two motivating scenarios along with relevant screen snapshots. Section 4 describes the implemen-

```

1 val log = "s3n://xcr:wJY@ws/logs/enroll.log"
2 val text_file = spark.textFile(log)
3 val avg = text_file
4   .map(line => (line.split()[2] ,
5               line.split()[3].toInt) )
6   .groupByKey()
7   .map(v => (v._1 , average(v._2)) )
8   .collect()

```

Figure 1: College student data analysis program in Scala

tation details of BIGDEBUG. Section 5 describes related work and Section 6 concludes. The current version of BIGDEBUG only supports Spark programs written in Scala and it is publicly available at <https://sites.google.com/site/sparkbigdebug/>.

2. BACKGROUND: APACHE SPARK

Apache Spark [2] is a large scale data processing platform that achieves orders-of-magnitude better performance than Hadoop MapReduce [1] for iterative workloads. BIGDEBUG targets Spark because of its wide adoption and support for interactive ad-hoc analytics. The Spark programming model can be viewed as an extension to the Map Reduce model with direct support for traditional relational algebra operators (e.g., group-by, join, filter). Spark programmers leverage Resilient Distributed Datasets (RDDs) to apply a series of transformations to a collection of data records (or tuples) stored in a distributed fashion e.g., in HDFS [11].

Calling a transformation on an RDD produces a *new* RDD that represents the result of applying the given transformation to the input RDD. Transformations are lazily evaluated. The actual evaluation of an RDD occurs, when an action such as `count` or `collect` is called. The Spark platform consists of three main entities: a *driver program*, a *master* node, and a set of *workers*. The master node controls distributed job execution and provides a rendezvous point between the driver and the workers. Internally, the Spark master translates a series of RDD transformations into a Directed Acyclic Graph (DAG) of *stages*, where each stage contains some sub-series of transformations, until a *shuffle step* is required (i.e., data must be re-partitioned). The Spark scheduler is responsible for executing each stage in topological order, with *tasks* that perform the work of a stage on input partitions. Each stage is fully executed before downstream dependent stages are scheduled. The final output stage evaluates the action that triggered the execution. The action result values are collected from each task and returned (via the master) to the driver program, which can initiate another series of transformations ending with an action.

3. MOTIVATING SCENARIOS WITH DEBUG FEATURES

Suppose Alice is a Spark user and she wants to process all US college student data. Because of the massive size of the data, she cannot store and analyze the data in a single machine. Suppose that she intends to compute the average age of all college students in each year (freshman, sophomore, junior, and senior). She starts by parsing the data into appropriate data types and then groups the records for each category. Once she has all related records grouped together, she computes the average and then collects the final results. A sample input record is in the following format:

```
1 Timothy Sophomore 21
```

The final program that Alice has written is shown in Figure 1. At line 2, she loads the US college student data from an Amazon S3 storage in the cluster. Line 4 reads each data record in the input data and generates a key value pair, where a key is the status category

for a student and the value is the age of that student. At lines 5 and 6, she groups all records with respect to the key and calculates the average for each category. At line 7, she executes the job and requests the result to be sent to the driver.

Simulated Breakpoint and Guarded Watchpoint

To maximize the throughput in a big data debugging session, BIGDEBUG provides **simulated breakpoints** that enable a user to inspect a program state in a remote executor node without actually pausing the entire computation. When such breakpoint is in place, a program state is regenerated, on-demand, from the last materialization point, while the original process is still running in the background. The last materialization point refers to the last stage boundary before the simulated breakpoint. These materialization points are determined beforehand by Spark’s scheduler.

To reduce developer burden in inspecting a large amount of intermediate records at a particular breakpoint within the workflow, BIGDEBUG’s **on-demand guarded watchpoints** retrieve intermediate data matching a user-defined predicate and transfer the selected data on demand. Furthermore, BIGDEBUG enables the user to update the guard predicate at runtime, while the job is still running. This dynamic guard update feature is useful when the user is not familiar with the data initially, and she wants to gradually narrow down the scope of the intermediate records to be inspected.

For example, suppose that Alice wants to inspect the program state at line 3. She can insert a simulated breakpoint using BIGDEBUG’s API i.e., `simulatedBreakpoint(r => !COLLEGEYEAR.contains(r.split()[2]))` with the guard predicate indicating that the second field is not one of the pre-defined college years. The benefit of this breakpoint combined with the guarded watchpoint is twofold. First, Alice can now inspect intermediate program results distributed across multiple nodes on the cloud, which is impossible in the original Spark. Second, she can also inspect records matching the guard predicate only, which tremendously reduces the inspection overhead.

While the Spark program instrumented with breakpoints is running on the cluster, Alice can use a web-based debugger interface by connecting to a configured port (4040 by default). Using this interface, she can view the DAG of the data flow program. On the left hand side of Figure 2, a yellow node in the DAG represents a breakpoint. Alice can use the code editor window on the right hand side to see the Spark program in execution. Statements with a breakpoint are tagged using a red arrow while the statement to be executed next is highlighted in blue.

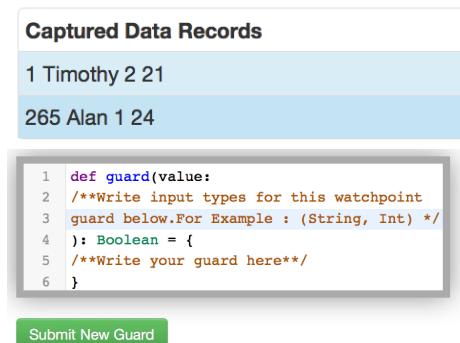


Figure 3: A user can edit the guard predicate using an editor.

Realtime Code Fix

After inspecting a program state at a breakpoint, if a user decides to patch code appearing later in the pipeline, she can use the **realtime code fix** feature to repair code on the fly. In this case, the original

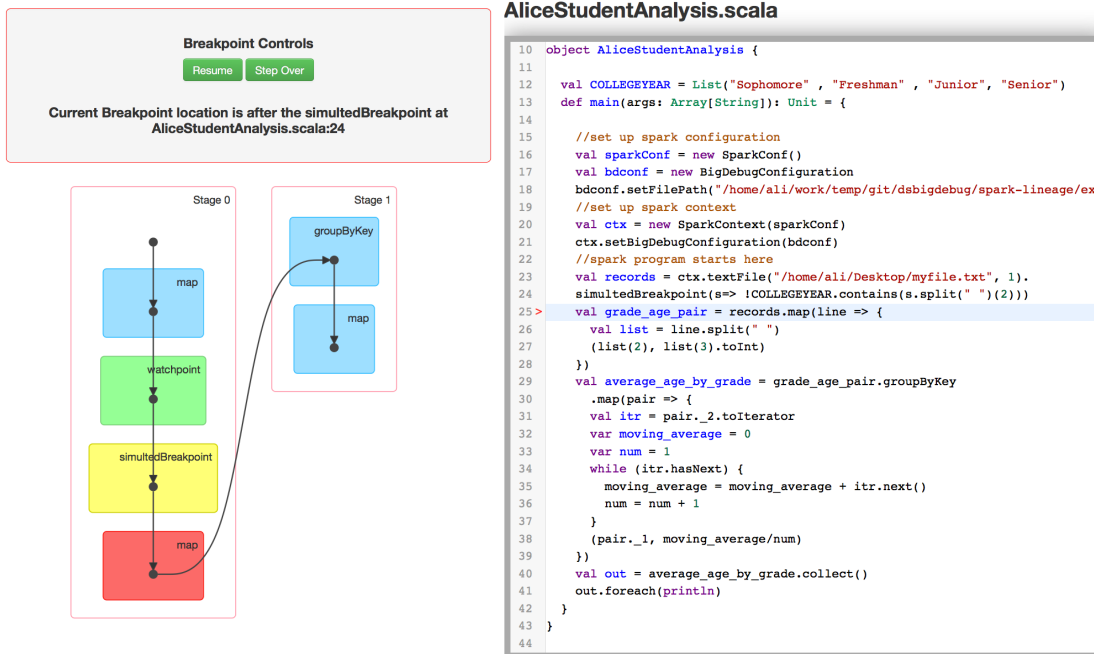


Figure 2: BIGDEBUG extends Spark’s user interface to provide runtime debugging features

job is canceled and a new job is created from the last materialization point before the breakpoint. This approach avoids restarting the entire job from scratch.

For example, in Figure 2, since a simulated breakpoint is in place, BIGDEBUG records the last materialization point before the breakpoint, in this case, after `textFile`. While the job is still executing, Alice can inspect the internal program state at the breakpoint. She can click on the green node on the DAG, which redirects her to a new web page, where intermediate records are displayed. When she requests to view the internal program state, the captured records from the guarded watchpoint are transferred to the driver node and displayed as shown in Figure 3. Upon viewing the intermediate records at a breakpoint, Alice discovers that some records use number 2 instead of Sophomore to indicate the status year.

```
1 Timothy 2 21
```

From this outlier record, Alice immediately learns that her program should handle records with a status year written in numbers. To apply realtime code fix, Alice can click on the corresponding transformation, in this case, `map` transformation marked in blue in the DAG. She can then insert a new user-defined function to replace the old user-defined function using a code editor shown in Figure 3. The code fix can now handle status year both in number and string formats. When Alice presses a submit button, BIGDEBUG compiles and redistributes the new user-defined function to each worker node and restarts the job from the latest materialization point. When the job finishes its execution, the final result after the fix is shown to Alice. In addition to a realtime code fix feature, Alice may use `resume` and `step over` as control commands. These control commands are available in BIGDEBUG’s UI.

Crash Culprit Remediation

In normal Spark, a runtime exception terminates the whole job, throwing away hours of computation while giving no information of the root cause of the error. When a Spark program fails with a runtime exception on the cluster, BIGDEBUG reports a **crash culprit** record in the intermediate stage but also identifies a **crash-inducing input(s)** in the original input data. While waiting for

a user intervention, BIGDEBUG runs pending tasks continuously to utilize idle resources in order to achieve high throughput. If a crash occurs, the original job keeps on running, while the user is notified of the fine-grained details of the crash. Once the crash culprit is reported to the user, the user can choose among three **crash remediation options**. First, a user can choose to **skip** the crash inducing record. The final output, in this case, will not reflect the skipped records. Second, a user can **modify crash culprit records** in realtime, so that the modified record can be injected back into the pipeline. Third, a user can **repair code**. The whole process of modifying crash culprits is optimized through lazy remediation. While the user takes time to resolve crash culprits, BIGDEBUG continues processing the rest of the records, while also reporting any additional crashing record. More details about crash remediation methods are discussed elsewhere [7].

Suppose that, after several hours of computation, a runtime exception occurs during the data processing. BIGDEBUG alerts Alice on the intermediate record responsible for the crash. These alerts turn the corresponding transformation node of the DAG to be red and highlight the corresponding code line in the main editor window to be red as well. When Alice clicks on the red node in the DAG, she is redirected to the crash culprit page of Figure 4. This page contains precise and useful information about the following crash culprit record:

```
1221 Matthew 4 24yr
```

When Alice is informed of the crash culprit record, BIGDEBUG continues executing the rest of the records and waits for the crash resolution from Alice. Alice may skip or modify the crash inducing intermediate record directly. Figure 4 shows the options provided on the UI to perform these remediation operations on the crash-inducing records. Alice skips the crashing record by pressing the `Skip` button on the crash culprit UI. BIGDEBUG also allows the batch repair of modifying all crash-inducing records at once using a user-defined repair script.

Forward and Backward Tracing

BIGDEBUG supports **fine-grained tracing** of individual records

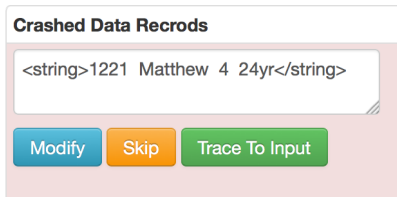


Figure 4: A user can either modify or skip the crash inducing records

by invoking a data provenance query on the fly. The *data provenance* problem in the database community refers to identifying the origin of final (or intermediate) output. Data provenance support for DISC systems is challenging, because operators such as aggregation, join, and `group-by` create many-to-one or many-to-many mappings for inputs and outputs and these mappings are physically distributed across different worker nodes. BIGDEBUG uses data provenance capability implemented through an extension of Spark’s RDD abstraction [8].

Fine-grained tracing allows users to reason about the faults in the program output or intermediate results, and explain why a certain problem has occurred. Using backward tracing, a crash culprit record can be traced back to the original inputs responsible for the crash record. Forward tracing allows user to find the output records affected by a selected input.

For example, during crash remediation, Alice can invoke forward and backward tracing feature at runtime to find the original input records responsible for the crash. On the crash culprit UI, Alice can invoke the backward tracing query by pressing the *trace to input* button. BIGDEBUG performs backward tracing in a new process to trace crash-inducing records in the original input data. Alice can also perform step-by-step backward tracing, showing all intermediate records tracing back to crash-inducing input records.

Fine-Grained Latency Monitoring

In big data processing, it is important to identify which records are causing delay. To localize performance anomalies at the record level, BIGDEBUG wraps each operator with a latency monitor. For each record at each transformation, BIGDEBUG computes the time taken to process each record, keeps track of a moving average, and sends a report to the monitor, if the time is greater than k standard deviations above the moving average, where default k is 2.

4. IMPLEMENTATION

The API for BIGDEBUG is shown in Figure 5 and targets Scala programs. All the features in BIGDEBUG is supported through the corresponding web-based user interface. BIGDEBUG extends the current Spark UI and provides a live stream of debugging information in an interactive and user-friendly manner. A screenshot of this interface is shown in Figure 2. Instead of creating a wrapper of existing Spark modules to track the input and output of each stage, BIGDEBUG directly extends Spark to monitor pipelined intra-stage transformations; its API extends the current RDD interface of Spark. A user can use function calls like `watchpoint()` and `simulatedBreakpoint()` on an RDD object to insert watchpoints and breakpoints. BIGDEBUG allows user to enable crash and latency monitoring on individual RDDs by calling appropriate methods on that RDD object. Tracing works at the granularity of each job and can be enabled or disabled through a `LineageContext`. All debugger control commands are linked with a driver that broadcasts the debugger control information to each worker. The runtime code patching is received and compiled at a driver and is then loaded into each worker, where an instru-

```

1 //RDD.scala
2 abstract class RDD[T: ClassTag] (
3   ...
4   def watchpoint(f: T => Boolean): RDD[T]
5   def simulatedBreakpoint
6   def simulatedBreakpoint(f:T => Boolean)
7   def enableLatencyAlert(set : Boolean)
8   def setCrashConfiguration(set :
9     CrashConfiguration)
9   def setFunction(f : T => U)
10  def goBackAll: LineageRDD
11  def goNextAll: LineageRDD
12  def goBack: LineageRDD
13  def goNext: LineageRDD
14  ...

```

Figure 5: BIGDEBUG’s API

mented task is running.

5. RELATED WORK

Fisher et al. interviewed 16 data analysts at Microsoft and studied the painpoints of big data analytics tools [5]. Their study finds that a cloud-based computing solution makes it far more difficult to debug. Xu et al. parse console logs and combine source code analysis to detect abnormal behavior [12]. Fu et al. map free-form text messages in log files to logging statements in source code [6]. None of these post-mortem log analysis approaches help developers debug DISC applications in realtime.

Inspector Gadget [9] is a framework proposal for monitoring and debugging data flow programs in Apache Pig [10]. The proposal is based on informal interviews with 10 Yahoo employees who write DISC applications. While Inspector Gadget proposes features such as step-through debugging, crash culprit determination, tracing, etc., it simply lists desired debug APIs, but leaves it to others to implement the proposed APIs. Arthur is a post-hoc instrumentation debugger that targets Spark and enables a user to selectively replay a part of the original execution [3]. However, a user can only perform *post-mortem* analysis and cannot inspect intermediate results at runtime. It also requires a user to write a custom query for post-hoc instrumentation. To localize faults, Arthur requires more than one run. In our prior work, we describe the design and evaluation of interactive debugging primitives [7] and data provenance for Apache Spark [8]. This demonstration paper builds on these prior works to showcase the UI and tool features of BIGDEBUG.

6. FUTURE WORK

BIGDEBUG offers interactive debugging primitives for an in-memory data-intensive scalable computing (DISC) framework. In terms of future work, instead of having a user specify a guard for an on-demand watchpoint, extracting data invariants from intercepted intermediate results may be useful for helping the user debug a DISC program. Another area for future work is tool-assisted automated fault localization in BIGDEBUG. For example, with the help of automated fault localization, we envision that a user can isolate the trace of failure-inducing workflow, diagnose the root cause of an error, and explain the cause-effect chain for unexpected results.

Acknowledgment

We thank the anonymous reviewers for their comments. Participants in this project are in part supported through NSF CCF- 1527923, CNS-1239498, CCF-1460325, CNS-1161595, IIS-1302698, CNS-1351047, and NIH U01HG008488.

7. REFERENCES

- [1] Hadoop. <http://hadoop.apache.org/>.
- [2] Spark. <https://spark.apache.org/>.
- [3] A. Dave, M. Zaharia, and I. Stoica. Arthur: Rich post-facto debugging for production analytics applications. Technical report, Citeseer, 2013.
- [4] J. Dean and S. Ghemawat. Mapreduce: simplified data processing on large clusters. *Communications of the ACM*, 51(1):107–113, 2008.
- [5] D. Fisher, R. DeLine, M. Czerwinski, and S. Drucker. Interactions with big data analytics. *interactions*, 19(3):50–59, May 2012.
- [6] Q. Fu, J.-G. Lou, Y. Wang, and J. Li. Execution anomaly detection in distributed systems through unstructured log analysis. In *Proceedings of the 2009 Ninth IEEE International Conference on Data Mining, ICDM '09*, pages 149–158, Washington, DC, USA, 2009. IEEE Computer Society.
- [7] M. A. Gulzar, M. Interlandi, S. Yoo, S. D. Tetali, T. Condie, T. Millstein, and M. Kim. Bigdebug: Debugging primitives for interactive big data processing in spark. In *Proceedings of the 38th International Conference on Software Engineering, ICSE '16*, pages 784–795, New York, NY, USA, 2016. ACM.
- [8] M. Interlandi, K. Shah, S. D. Tetali, M. A. Gulzar, S. Yoo, M. Kim, T. D. Millstein, and T. Condie. Titian: Data provenance support in spark. *PVLDB*, 9(3):216–227, 2015.
- [9] C. Olston and B. Reed. Inspector gadget: A framework for custom monitoring and debugging of distributed dataflows. In *Proceedings of the 2011 ACM SIGMOD International Conference on Management of data*, pages 1221–1224. ACM, 2011.
- [10] C. Olston, B. Reed, U. Srivastava, R. Kumar, and A. Tomkins. Pig latin: a not-so-foreign language for data processing. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, pages 1099–1110. ACM, 2008.
- [11] J. Shafer, S. Rixner, and A. L. Cox. The hadoop distributed filesystem: Balancing portability and performance. In *Performance Analysis of Systems & Software (ISPASS), 2010 IEEE International Symposium on*, pages 122–133. IEEE, 2010.
- [12] W. Xu, L. Huang, A. Fox, D. Patterson, and M. I. Jordan. Detecting large-scale system problems by mining console logs. In *Proceedings of the ACM SIGOPS 22nd symposium on Operating systems principles*, pages 117–132. ACM, 2009.
- [13] M. Zaharia, M. Chowdhury, T. Das, A. Dave, J. Ma, M. McCauley, M. J. Franklin, S. Shenker, and I. Stoica. 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*, pages 2–2. USENIX Association, 2012.