White-Box Testing of Big Data Analytics with Complex User-Defined Functions

Muhammad Ali Gulzar ¹

Shaghayegh Mardani¹ Madan Musuvathi² Miryung Kim¹

¹University of California, Los Angeles ²Mircrosoft Research



Inadequate Testing of Big Data Analytics

Software Development Cycle of Big Data Analytics



Motivating Example

Find the total number of trips made from UCLA using a public transport, a personal vehicle, or on foot.

Trips Dataset (20GB)							
<u></u> #,	<u>ORIG</u> ,	<u>DEST</u> ,	<u>DIST</u> ,	TIME			
1,	90034,	90024,	10,	1			
2,	90001,	90024,	16,	1.4			

Locations Dataset (100MB)

<u>Zip</u> ,	<u>Location</u>
90034,	"UCLA"
90024,	"Westwood"

Big Data Application in Apache Spark

...

Characteristics of Big Data Analytics



How do we test a big data application effectively and efficiently?

Option 1: Sample Input Data

- random sampling,
- top n sampling
- top k% sample, etc.

Limitations:

- The sample may only exercise a limited set of **program paths** (low code coverage).
- The sample may not include the inputs leading to a **program crash**.
- A large sample may have higher coverage but increase local **testing time**.



Option 2: Traditional Test Generation for Java

- Big Data Analytics programs compile to Java bytecode
- But this includes the entire system (700 KLOC for Apache Spark)



• Symbolic execution without abstraction is infeasible and would not scale

Our Approach: White-Box Testing



- 1. Decompose relational skeleton and UDFs
- 2. Logical specifications for relational operators
- 3. Symbolic execution of UDFs
- 4. Generate inputs by joint path constraints

Step 1Step 2:StaDecompositionLogical SpecsEven

olic Ster Ster

Modelling Dataflow Operators



p 1 Step position Logica Step 4: Test Generation

Modelling Dataflow Operators



- Handle terminating and nonterminating cases of dataflow operators
- E.g. Join can introduce 3 cases
 - 2 cases in which keys from right and left do not match
 - 1 case in which right and left keys match

Step 1Step 2:Step 3:Step 4: TestDecompositionLogical SpecsExecutionGeneration

Modelling User-defined Functions



Step 1 ecomposition ep 2: S al Specs Step 4: Test Generation

Join Dataflow and UDF (JDU) Path





Test Input Generation



Evaluation

RQ1: How much test coverage improvement can BigTest achieve?

RQ2: How many faults can BigTest detect?

 We built the first benchmark of faulty dataflow programs based on our survey of such programs on Q/A forums *e.g.* StackOverflow .

RQ3: How much test data reduction does BigTest provide and how long does BigTest take to generate test data?

Experimental Setting

- We use seven subject programs from earlier works
- All subject applications have complex string, complex arithmetic, Tuple type for key-value pairs, and collections with custom logic.

Subject Program	Dataflow Operators	# of Operators	JDU Paths K=2	# of UDFs
Income Aggregate	map, filter, reduce	3	6	4
Movie Ratings	map, filter, reduceByKey	4	5	4
Airport Layover	map, filter, reduceByKey	3	14	4
Commute Type	map, fitler, join, reduceByKey	6	11	5
PigMix-L2	map, join	5	4	6
Grade Analysis	flatmap, filter, reduceByKey, map	5	30	3
Word Count	flatmap, map, reduceByKey	3	4	3

Study of Big Data Analytics Faults



- No existing benchmark of faulty applications
- We study the characteristics of real-world big data analytics bugs posted on **StackOverflow** and **Apache Spark Mailing Lists**.

Survey Statistics			Fault Types	Example		
Keywords Searched	Apache Spark exceptions,		Incorrect String Offset	<pre>str.substring(1,0)</pre>		
	task errors, failures, wrong outputs		Incorrect Column Selection	<pre>str.split(",")[1]</pre>		
Posts Studied	Тор 50		Wrong Delimiters	<pre>str.split("\t")[1]</pre>		
Posts with Coding	23		Incorrect Branch Condition	If(age>10 && age<9)		
Errors			Wrong Join Type	LeftOuterJoin		
Common Fault Types	7		Key-Value Swap	(Value, Key)		
Total Faulty Programs 31			Others	Division by zero		

Real World Fault Injection

- Identified 7 common code fault types
- Manually inserted these faults into benchmarks
- Leads to a total of **31 faulty big data applications**.

Original Program

```
val trips = sc.textFile("trips")
.map { s =>
    val c = s.split(",");
    (c(1), c(3).toInt / c(4).toInt)
}
val loc = sc.textFile("zipcode")
....
```

Faulty Program

val trips = sc.textFile("trips") .map { s => val c = s.split(","); (c(1), c(2).toInt / c(5).toInt) } val loc = sc.textFile("zipcode")

After injecting fault based on fault type *"Incorrect Column Selection"*, the program extracts the column at index 5 instead of 4.

RQ1: Code Coverage



Sedge [ASE '13] generates examples for dataflow programs but it handles a UDF as **uninterpreted** function and **does not model its internals**.

RQ1: Code Coverage



BigTest improves JDU path coverage by 78% against Sedge and 34% against the entire dataset.

RQ2

RQ2: Fault Detection Capability

Applications	Total Detected by Detected					Injected Fault Type					
	Seeded Faults	BigTest	by Sedge	1	2	3	4	5	6	7	
Income Aggregate	3	3	1	1	NA	NA	\checkmark	NA	NA	1	
Movie Rating	6	6	6	\checkmark	\checkmark	\checkmark	\checkmark	NA	\checkmark	\checkmark	
Airport Layover	6	6	4	1	✓	\checkmark	✓	NA	\checkmark	\checkmark	
Commute Type	6	6	4	NA	\checkmark	1	\checkmark	1	\checkmark	\checkmark	
PigMix-L2	4	4	2	NA	\checkmark	\checkmark	NA	\checkmark	\checkmark	NA	
Grade Analysis	4	4	3	NA	\checkmark	\checkmark	\checkmark	NA	NA	\checkmark	
Word Count	2	2	0	NA	1	NA	NA	NA	NA	1	

BigTest detects 2X more faults than Sedge because it models the internal semantics of UDFs with the specifications of dataflow operators.

RQ1 RQ2 RQ3

RQ3: Test Size Reduction



Compared to the entire dataset, BigTest achieves more JDU path coverage with 10⁵X to 10⁸X smaller test data, translating in to 194X testing speed up.

Summary

- Need SE tools for big data analytics applications
- BigTest provides *exhaustive*, *automatic*, and *fast* testing
- Contributions:
 - 1. Demonstrated the need to interpret UDFs
 - 2. Model strings, collections, and tuples
 - 3. Logical specifications for dataflow operators handling terminating and nonterminating cases
 - 4. Provide the first symbolic execution engine for Apache Spark/Scala
 - 5. Present a study of big data analytics bugs and the first bug benchmark

Publically available at: https://github.com/maligulzar/BigTest

RQ3: Breakdown of BigTest's Testing Time

Breakdown of Testing Time



By running tests locally, BigTest improves the testing time (CPU seconds) by 194X, on average, compared to testing the entire dataset on 16-node cluster.

Inadequate Test Generation Tools for Big Data Analytics

Traditional Software Test Generation

```
def concat(append: boolean, a:
String, b: String ) {
   result: String = null;
   If (append)result = a + b;
   return
result.toLowerCase();
}
```

- Standalone application
- Symbolic Execution Compatible
- Well defined semantics
- Logical execution is similar to physical execution

Big Data Analytics Test Generation

sc.textFile("hdfs")
.flatMap(s=> s.split(","))
.map(w =>(w,1))
.reducebyKey(_+_)

- Heavily depends on framework
- Non-existence Symbolic Execution for dataflow operators
- New operators with changing semantics
- Logical execution is different to physical execution

24

Program Decomposition

- **Challenge:** Due to the complexity of DISC frameworks' code, symbolic execution is infeasible on DISC applications.
- Insight: The individual UDFs of DISC application are relatively smaller (<100 LOC) making symbolic execution feasible.
- **Solution:** We decompose a DISC application using AST analysis into a set of individual UDFs and dataflow operators.



Symbolic Execution of UDFs

- **Challenges:** Strings, Collections, and Object are eminent in DISC applications but not fully support by symbolic execution tool *i.e* Java PathFinder.
- Insight: In DISC applications, most unbounded types are eventually bounded. We perform lazy SE on such types *e.g* Split(",") is unbounded Array but Split(",")[1] is bounded.
- **Solution:** Using JPF, we symbolically execute UDFs in isolation to generated path constraints and effects. Loops and Arrays are bounded by K=2.



Logical Specifications of Dataflow Operators

- **Challenges:** Dataflow operators in DISC applications are accompanied with 100Ks lines of framework code making symbolic execution infeasible.
- **Insight:** Dataflow operators have standard semantics but implemented differently for optimization purposes.
- **Solution:** Using these semantics, we abstract their implementation in logical specifications and used the specifications to tie together UDFs' symbolic trees.

