SparkFuzz: Searching Correctness Regressions in Modern Query Engines

Bogdan Ghit, **Nicolas Poggi**, Josh Rosen, Reynold Xin, and Peter Boncz^{*}

June 19 - DBTest 2020



databricks UNIFIED DATA ANALYTICS PLATFORM



DATA SCIENCE WORKSPACE



ENTERPRISE CLOUD SERVICE



Introduction



Apache Spark

Fast and expressive data processing engine

- distributed computing
- rich APIs
 - including SQL
- large community

Started at UC Berkeley in 2009

- 2010 open sourced
- 2014 top level project
- 2020 v3 released (10 years!)

June 2002 v 3.0.0 released

3500+ resolved tickets



SparkFuzz proposal

1. Leverage fuzz testing techniques a. to complement SQL testing b. automate bug discovery Design of a toolkit for SQL engines 2. a. model for randomized i. DDL, data, and queries b. A runner and evaluator 3. Applicability of coverage metrics a. as test stop gaps reducing time (and costs) b. enabling more testing dimensions c.



DDL and data generation

Random number of columns

Automated dataset generation

- by randomly sampling
 - supported data types
 - parameter ranges
- Producing valid schemas
- . Populating datasets





Recursive query model w/ a probabilistic profile



Operators and features annotated with:

Independent weights

Optional clauses

Inter-dependent weights

- Join types
- Select functions

Query and regression example

Query produced in a small dataset with 2 tables of 5x5 size

- Within 10 queries, this query triggered an exception
- Related to COALESCE flattening

Correctness regression example [SPARK-16633]

Using constant input values breaks the the LEAD function

Spark [1.0, 696, -871.81, -64.98, -349]
PostgreSQL [1.0, 696, -871.81, NULL, -349]

Query operator coverage analysis



Continuous Integration pipeline





Conclusion and future work



- Prevented SQL correctness errors reaching production
 - complementing the testing practices
- Runtime operator coverage metrics found applicable
 - For testing code changes rapidly
 - With a degree of coverage
- Future work
 - Improve the metric coverage to include operator chaining
 - Update the model generation to use Spark AST grammar directly

SparkFuzz: Searching Correctness Regressions

Thanks, questions?

Bogdan Ghit, Nicolas Poggi, Josh Rosen, Reynold Xin, and Peter Boncz

Feedback: Nicolas.Poggi@databricks.com

latabricks