# Correctness and Performance of Apache Spark SQL

Spark + AI Summit, London

October 4, 2018

databricks

# About us

## BOGDAN GHIT

Databricks, Software Engineer
- SQL performance optimizations

IBM T.J. Watson, Research Intern
- Bid advisor for cloud spot markets

Delft University of Technology, PhD in Computer Science
- Resource management in datacenters
- Performance of Spark, Hadoop

## NICOLAS POGGI

Databricks, Performance Engineer
- Spark benchmarking

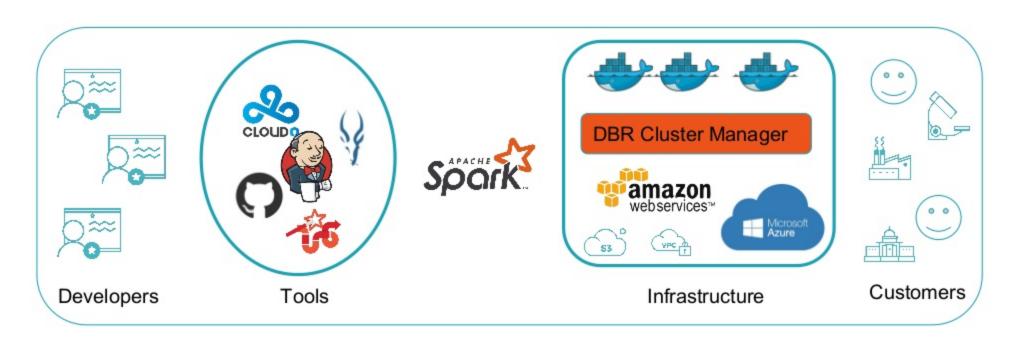Barcelona Supercomputing - Microsoft Research Centre
- Lead researcher ALOJA project
- New architectures for Big Data

BarcelonaTech (UPC), PhD in Computer Architecture
- Autonomic resource manager for the cloud
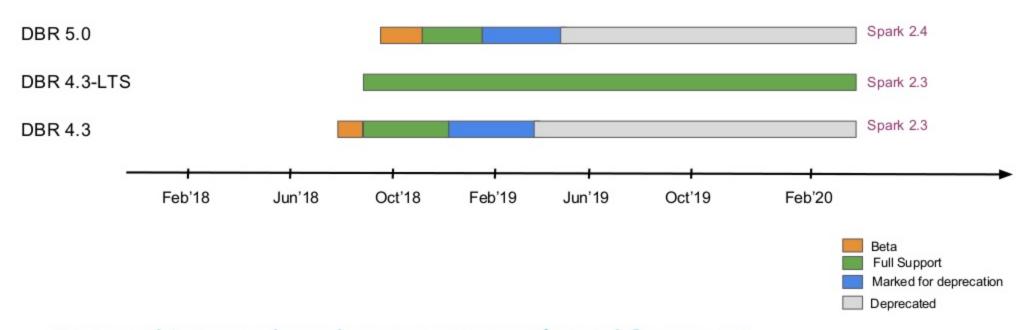- Web customer modeling

# Databricks ecosystem



Developers

Tools

DBR Cluster Manager

Infrastructure

Customers

# Databricks runtime (DBR) releases



DBR 5.0 — Spark 2.4
DBR 4.3-LTS — Spark 2.3
DBR 4.3 — Spark 2.3

Feb'18    Jun'18    Oct'18    Feb'19    Jun'19    Oct'19    Feb'20

Beta
Full Support
Marked for deprecation
Deprecated

## Our goal is to make releases **automatic** and **frequent**

databricks

# Apache Spark contributions



Number of commits

Hundreds of commits monthly to the Apache Spark project

At this pace of development, **mistakes** are bound to happen

# Where do these contributions go?

SQL Core

Over 200 built-in functions

Scope of the testing

Unit tests
50.3%

Source code
49.7%

Query

Input data

Configuration
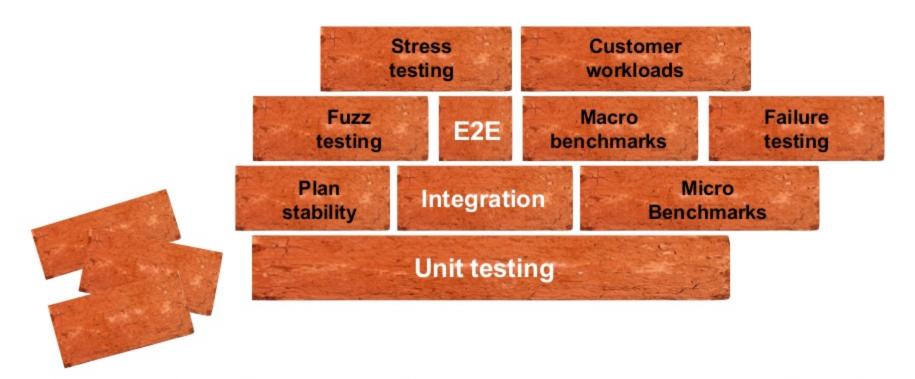
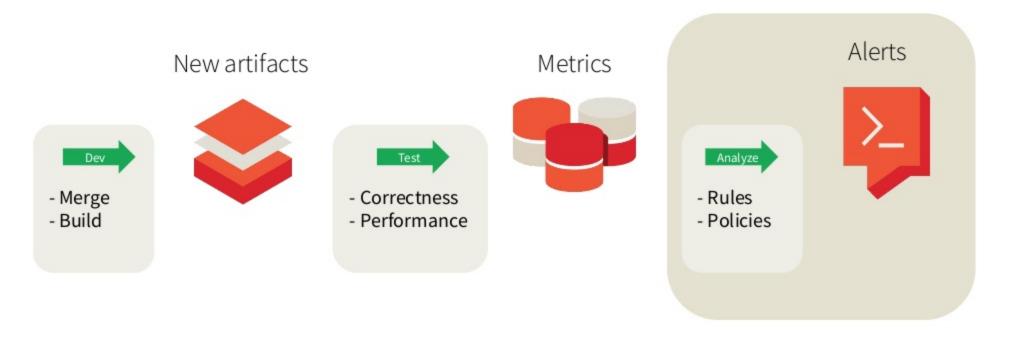Developers put a significant **engineering effort** in testing

databricks

6

# Yet another brick in the wall



Unit testing *is not enough* to guarantee correctness and performance

# Continuous Integration pipeline

New artifacts

Metrics

Alerts

Dev
- Merge
- Build

Test
- Correctness
- Performance

Analyze
- Rules
- Policies

>_

databricks

# Classification and alerting



Events → Classify

✓

- Impact
- Scope
- Correlation
- Confirm?

Failure → Correctness

Regression → Performance

Re-test → Root-cause

- Minimize
- Drill-down
- Profile
- Compare
- Validate

→ Alert

databricks

# Correctness

# Random query generation

# DDL and datagen

String

Boolean

Decimal

Integer

Timestamp

BigInt

SmallInt

Float

Choose a data type

Random number of columns

Random partition columns

Random number of rows

Random number of tables

# Recursive query model

Query

**SQL Query**

Clause

WITH

UNION

SELECT

WHERE

GROUP BY

ORDER BY

JOIN

FROM

Expression

Functions

Constant

Alias

Column

Table

databricks

13

# Probabilistic query profile

## Independent weights

- Optional query clauses

10% UNION   10% GROUP BY   50% WHERE   10% ORDER BY

## Inter-dependent weights

- Join types
- Select functions



FULL-OUTER 2.2%
RIGHT 7.5%
LEFT 22.4%
INNER 67.2%

databricks

# Coalesce flattening (1/4)

```sql
SELECT COALESCE(t2.smallint_col_3, t1.smallint_col_3, t2.smallint_col_3) AS int_col,
    IF(NULL, VARIANCE(COALESCE(t2.smallint_col_3, t1.smallint_col_3, t2.smallint_col_3)),
    COALESCE(t2.smallint_col_3, t1.smallint_col_3, t2.smallint_col_3)) AS int_col_1,
    STDDEV(t2.double_col_2) AS float_col,
    COALESCE(MIN((t1.smallint_col_3) - (COALESCE(t2.smallint_col_3, t1.smallint_col_3,
    t2.smallint_col_3))), COALESCE(t2.smallint_col_3, t1.smallint_col_3, t2.smallint_col_3),
    COALESCE(t2.smallint_col_3, t1.smallint_col_3, t2.smallint_col_3)) AS int_col_2
FROM table_4 t1
INNER JOIN table_4 t2 ON (t2.timestamp_col_7) = (t1.timestamp_col_7)
WHERE (t1.smallint_col_3) IN (CAST('0.04' AS DECIMAL(10,10)), t1.smallint_col_3)
GROUP BY COALESCE(t2.smallint_col_3, t1.smallint_col_3, t2.smallint_col_3)
```

Small dataset with 2 tables of 5x5 size
Within 10 randomly generated queries

Error: Operation is in ERROR_STATE

databricks

# Coalesce flattening (2/3)



Aggregate — `GROUP BY COALESCE(foo.id, foo.val)`

Project — `COALESCE(COALESCE(foo.id, foo.val), 88)`

Join — `foo.ts  = bar.ts`

SCAN foo

FILTER — `foo.id IN`
`(CAST('0.04' AS DECIMAL(10, 10)), foo.id)`

SCAN bar

databricks

# Coalesce flattening (3/4)

```
Aggregate          COALESCE(foo.id, foo.val)

  Project          COALESCE(COALESCE(foo.id, foo.val), 88)

  Join             foo.ts  = bar.ts

SCAN t1   FILTER   foo.id IN
                   (CAST('0.04' AS DECIMAL(10, 10)), foo.id)

          SCAN t2
```

# Coalesce flattening (4/4)

Aggregate

Project

SCAN foo

**Minimized query:**

```
SELECT
    COALESCE(COALESCE(foo.id, foo.val), 88)
FROM foo
GROUP BY
    COALESCE(foo.id, foo.val)
```

Analyzing the error
- The optimizer flattens the nested coalesce calls
- The SELECT clause doesn't contain the GROUP BY expression
- Possibly a problem with any GROUP BY expression that can be optimized

# Lead function (1/3)

```sql
SELECT (t1.decimal0803_col_3) / (t1.decimal0803_col_3) AS decimal_col,
       CAST(696 AS STRING) AS char_col, t1.decimal0803_col_3,
       (COALESCE(CAST('0.02' AS DECIMAL(10,10)),
                 CAST('0.47' AS DECIMAL(10,10)),
                 CAST('-0.53' AS DECIMAL(10,10)))) +
       (LEAD(-65, 4) OVER (ORDER BY (t1.decimal0803_col_3) / (t1.decimal0803_col_3),
                 CAST(696 AS STRING))) AS decimal_col_1,
                 CAST(-349 AS STRING) AS char_col_1

FROM table_16 t1
WHERE (943) > (889)
```
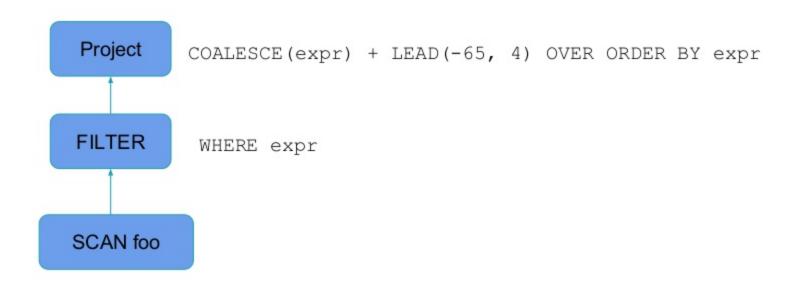
**Error:** Column 4 in row 10 does not match:

      [1.0, 696, -871.81, <<-64.98>>, -349] SPARK row

      [1.0, 696, -871.81, <<None>>, -349] POSTGRESQL row

# Lead function (2/3)



```
Project        COALESCE(expr) + LEAD(-65, 4) OVER ORDER BY expr

FILTER         WHERE expr

SCAN foo
```

# Lead function (3/3)

```
Project    COALESCE(expr) + LEAD(-65, 4) OVER ORDER BY expr

FILTER     WHERE expr

SCAN foo
```
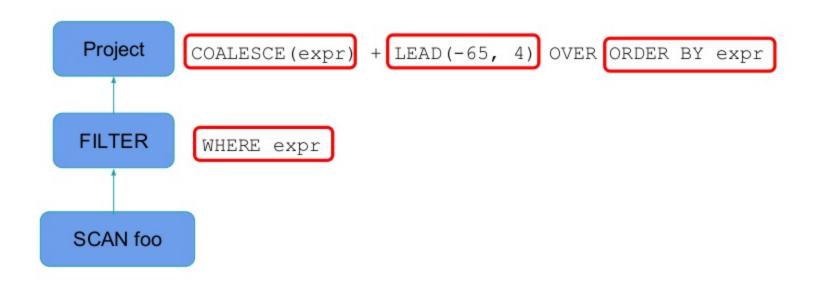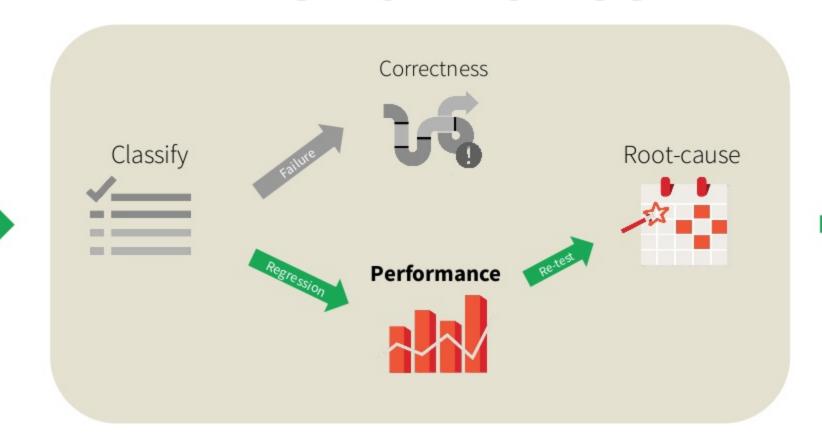
Analyzing the error
- Using constant input values breaks the behaviour of the LEAD function
- SC-16633: https://github.com/apache/spark/pull/14284

databricks

# Performance

# Benchmarking tools

- We use spark-sql-perf public library for TPC workloads
  - Provides datagen and import scripts
    - local, cluster, S3
  - Dashboards for analyzing results

- The Spark micro benchmarks
- And the async-profiler
  - to produce flamegraphs



### Spark SQL Performance Tests

build passing

This is a performance testing framework for Spark SQL in Apache Spark 2.2+.

**Note: This README is still under development. Please also check our source code for more information.**

### Quick Start

**Running from command line.**

https://github.com/databricks/spark-sql-perf



*CPU Flame Graph*

http://www.brendangregg.com/flamegraphs.html

# DBR 5.0-beta (v2.4) performance ~~tracking~~ -- **journey**



Dashboard showing:

- daysBack: 30 | daysTargetBack: 30 | targetVersion: master-2.4

**Time series of runs by type (10 days)** — Version: DBR 4.3 (2.3), master-2.3, master-2.4. Arrow marker with "15%" annotation near Aug 15.

**Top 5 regressing queries to previous version** — bar chart of GM% difference for queryName q78, q93, q50, q64, q84.

**Per run and type geomean** — quantile scatter plots for DBR 4.3 (2.3), master-2.3, master-2.4 showing Geomean (s) vs Quantile.

# Per query drill-down: 67

## First, **scope** and **validate**



Query 67: **18%** regression    From 320s to 390s

Query, version
- q67, 2.4-master
- q67, 2.3-master

- in 2.4-master (dev) compared
- to 2.3 in DBR 4.3 (prod)

# Q67 executor profile for Spark 2.4-master



org/apache/spark/shuffle/sort/BypassMergeSortShuffleWriter.write
org/apache/spark/scheduler/ShuffleMapTask.runTask
org/apache/spark/scheduler/ShuffleMapTask.runTask
org/apache/spark/scheduler/Task.run
org/apache/spark/executor/Executor$TaskRunner.run
java/util/concurrent/ThreadPoolExecutor.runWorker
java/util/concurrent/ThreadPoolExecutor$Worker.run
java/lang/Thread.run

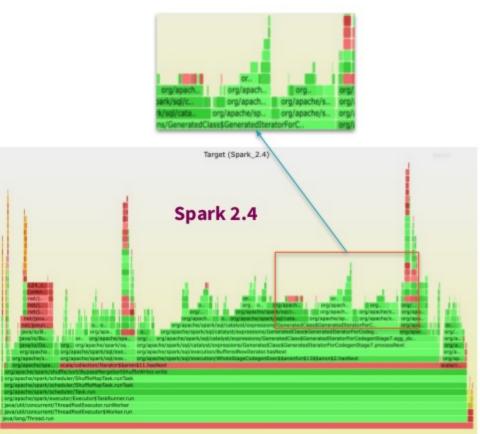# Side-by-side 2.3 vs 2.4: find the differences

# Framegraph diff zoom

**Look for hints:**

- Mem mgmt
- Hashing
- unsafe



Reset Zoom

Diff-After

Blue = speedup, Red == slowdown

Search

Murmur3_x86_32.hashUnsafeBytesBlock()

unsafe/BytesToBytesMap.safeLookup

**New:** hash/Murmur3_x86_32.hashUTF8String()

unsafe/Platform.copyMemory()

org/apache/spark/sql/catalyst/expressions/GeneratedClass$GeneratedIteratorForCodegenStage7.agg_doConsume_0$
org/apache/spark/sql/catalyst/expressions/GeneratedClass$GeneratedIteratorForCodegenStage7.expand_doConsume_0$
org/apache/spark/sql/catalyst/expressions/GeneratedClass$GeneratedIteratorForCodegenStage7.agg_doAggregateWithKeys_0$
org/apache/spark/sql/catalyst/expressions/GeneratedClass$GeneratedIteratorForCodegenStage7.processNext
org/apache/spark/sql/execution/BufferedRowIterator.hasNext
org/apache/spark/sql/execution/WholeStageCodegenExec$$anonfun$13$$anon$2.hasNext
scala/collection/Iterator$$anon$11.hasNext
org/apache/spark/shuffle/sort/BypassMergeSortShuffleWriter.write
org/apache/spark/scheduler/ShuffleMapTask.runTask
org/apache/spark/scheduler/ShuffleMapTask.runTask
org/apache/spark/scheduler/Task.run
org/apache/spark/executor/Executor$TaskRunner.run
java/util/concurrent/ThreadPoolExecutor.runWorker
java/util/concurrent/ThreadPoolExecutor$Worker.run
java/lang/Thread.run
all

databricks

# Root-causing

## GIT BISECT

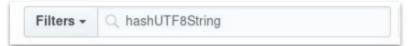Microbenchmark for UTF8String

```
test("hashing") {
    import org.apache.spark.unsafe.hash.Murmur3_x86_32
    import org.apache.spark.unsafe.types.UTF8String
    val hasher = new Murmur3_x86_32(0)
    val str = UTF8String.fromString("b" * 10001)
    val numIter = 100000
    val start = System.nanoTime
    for(i <- 0 until numIter) {
        Murmur3_x86_32.hashUTF8String(str, 0)
```
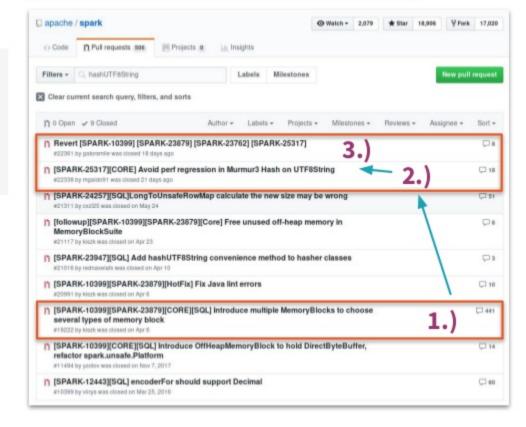
**Results:**

- Spark 2.3: hashUnsafeBytes() -> **40μs**
- Spark 2.4 *hashUnsafeBytesBlock()* -> **140μs**
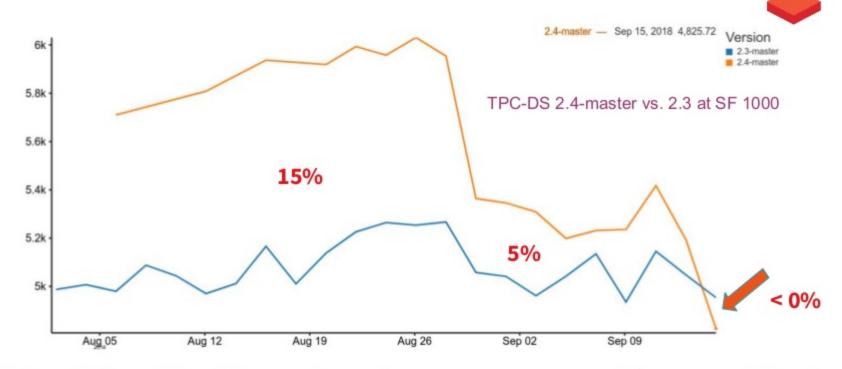
- also slower *UTF8String.getBytes()*



databricks

# It is a journey to get a release out



2.4-master — Sep 15, 2018 4,825.72

Version
- 2.3-master
- 2.4-master

TPC-DS 2.4-master vs. 2.3 at SF 1000

15%

5%

< 0%

**DBR and Spark testing and performance are a continuous effort**
- **Over a month** effort to bring performance to improving

databricks

# Conclusion

Spark in production is *not just the framework*

Unit and integration testing are not enough

We need Spark specific tools to automate the process

to ensure both correctness and performance

# Thanks!

Dev → Test → Analyze

## Correctness and Performance of Apache Spark SQL

October 2018

databricks