# Correctness and Performance of Apache Spark SQL

Spark + AI Summit, London

October 4, 2018

databricks®

# About us

**BOGDAN GHIT**

**NICOLAS POGGI**

Databricks, Software Engineer
- SQL performance optimizations

Databricks, Performance Engineer
- Spark benchmarking

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

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

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

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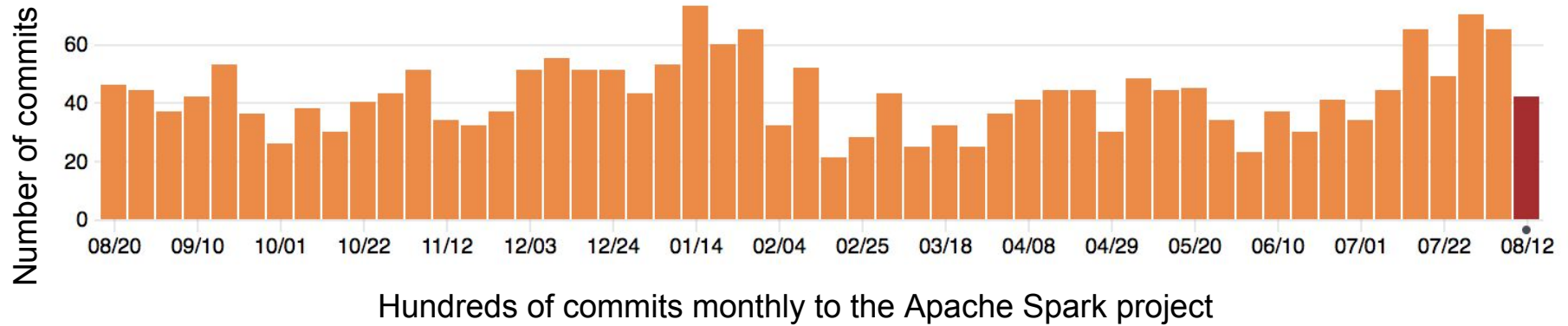- Web customer modeling

# 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

**Legend:**
- 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

databricks

# Where do these contributions go?

SQL Core



Unit tests
50.3%

Source code
49.7%

Scope of the testing

Over 200 built-in functions

Query
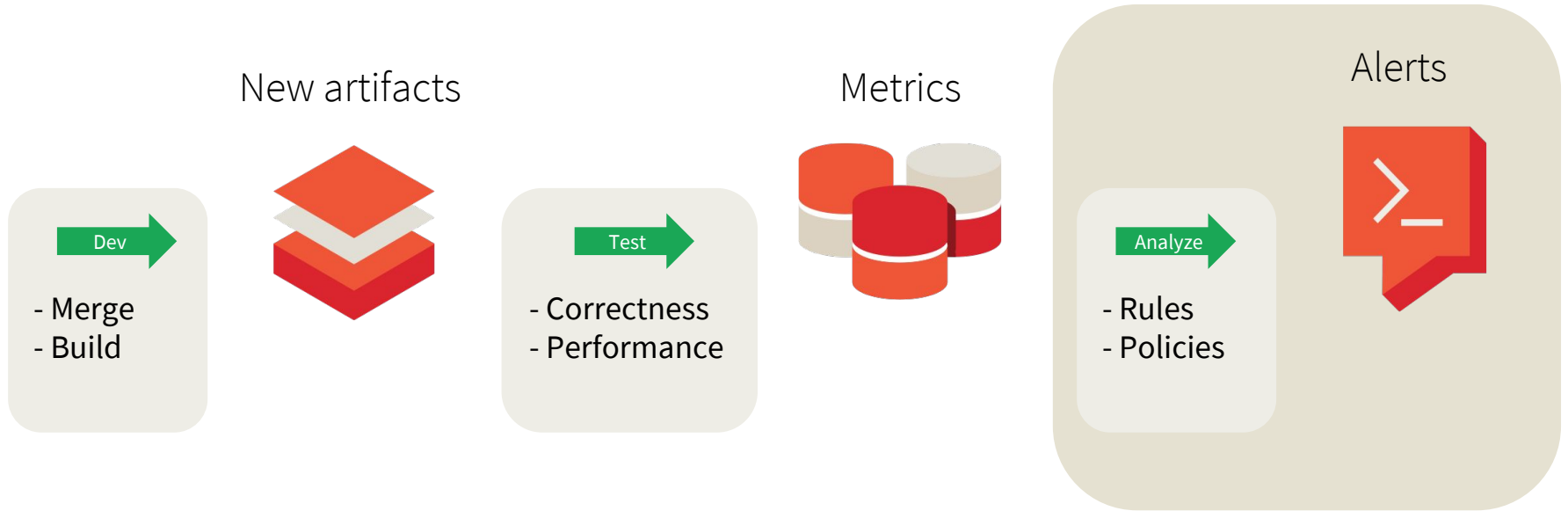
Input data

Configuration

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

databricks

6

# Yet another brick in the wall
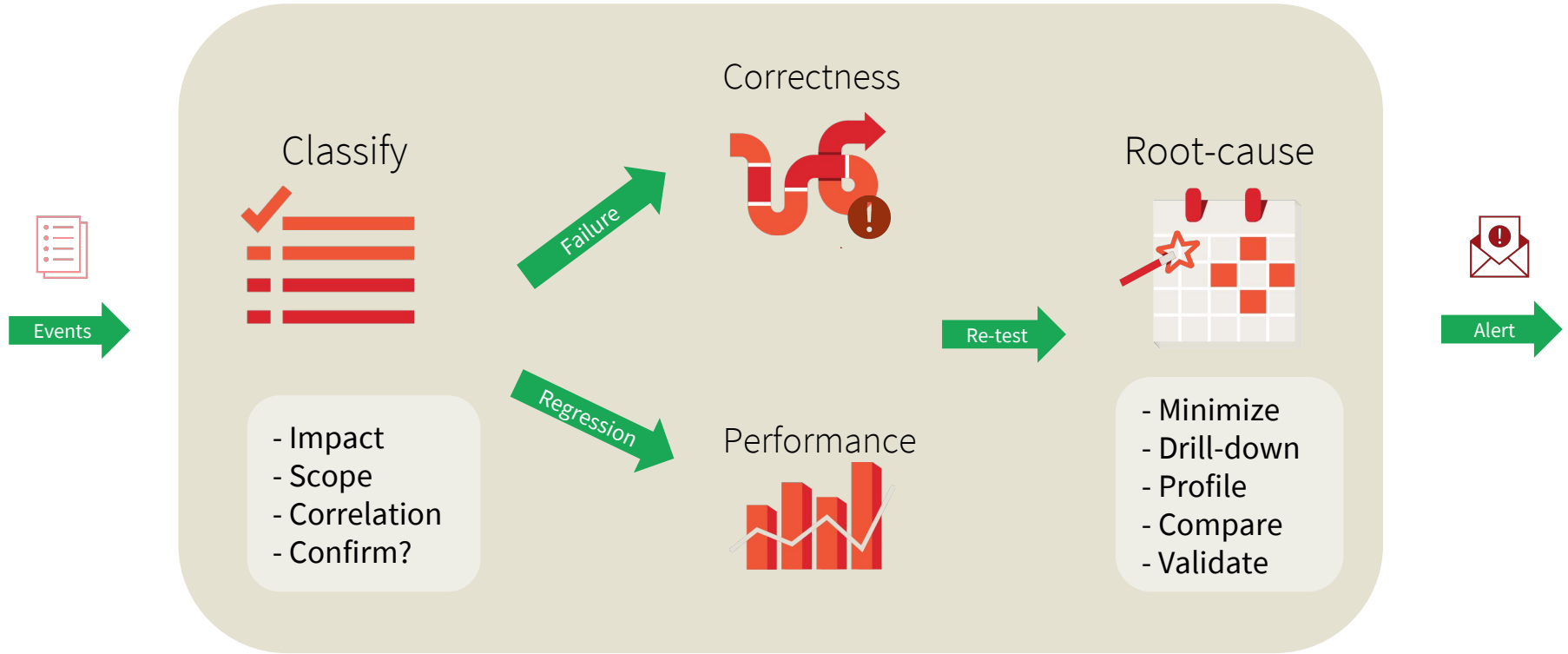


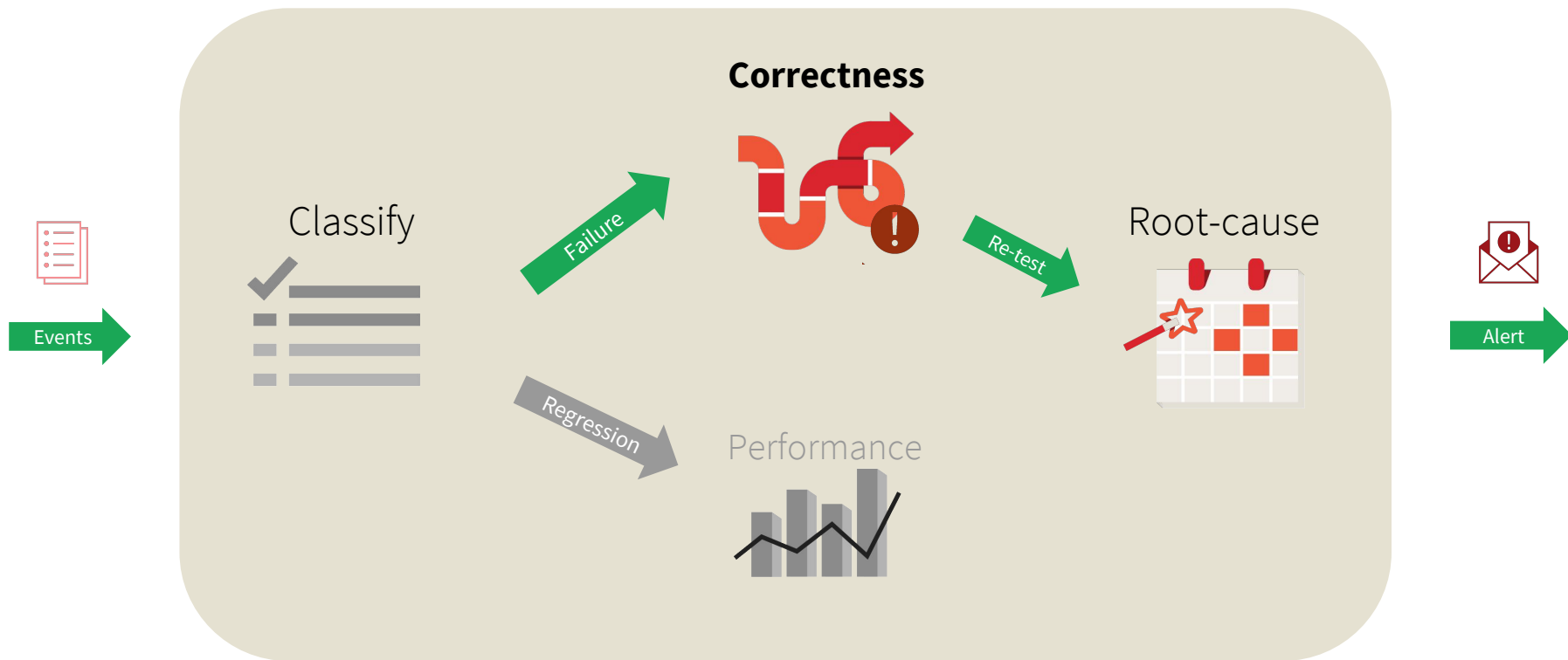Unit testing *is not enough* to guarantee correctness and performance

# Continuous Integration pipeline

New artifacts

Metrics

Alerts

Dev
- Merge
- Build

Test
- Correctness
- Performance

Analyze
- Rules
- Policies

databricks™

8

# 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™

9

# Correctness

# Random query generation

# DDL and datagen

Random number of columns

Random partition columns

Choose a data type

String

Boolean

BigInt

Decimal

SmallInt

Integer

Float

Timestamp

...

...

...

Random number of rows

Random number of tables

databricks

12

# Recursive query model

Query

SQL Query

Clause

WITH
UNION
SELECT
WHERE
GROUP BY
ORDER BY
JOIN
FROM

Expression

Functions
Constant
Alias
Column
Table

# Probabilistic query profile

## Independent weights

- Optional query clauses

10%
UNION

GROUP BY
10%

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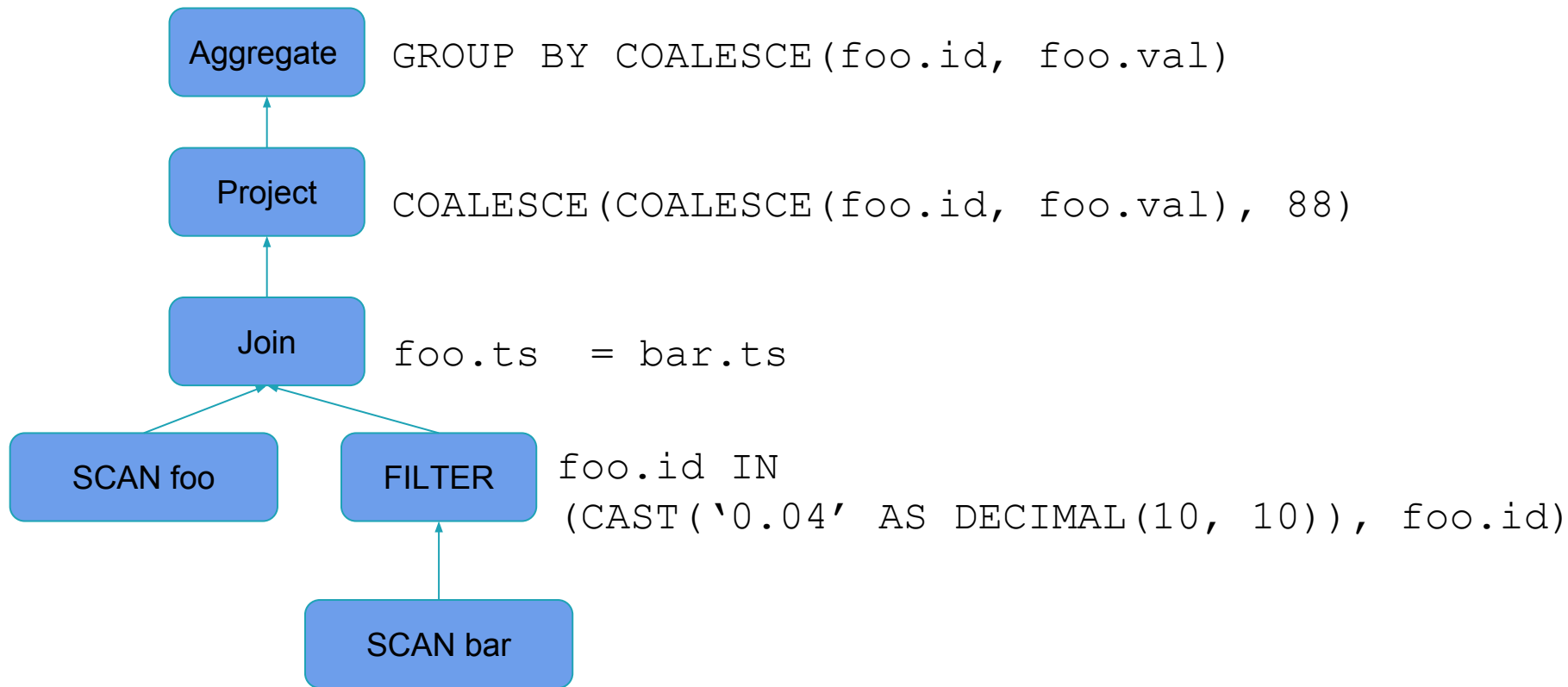
Small dataset with 2 tables of 5x5 size
Within 10 randomly generated queries

Error: Operation is in ERROR_STATE

# 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)



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
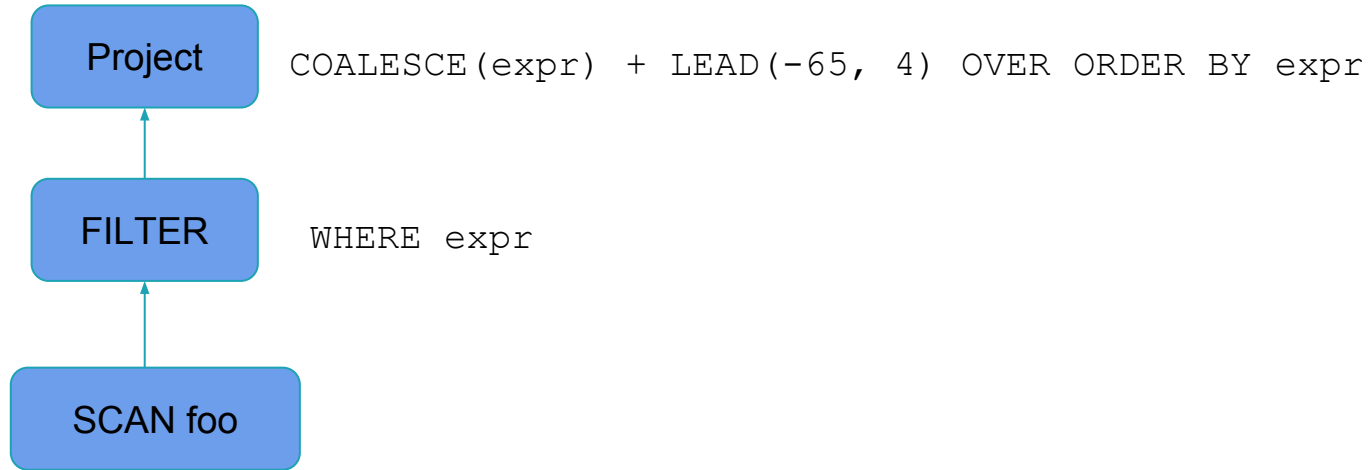
databricks

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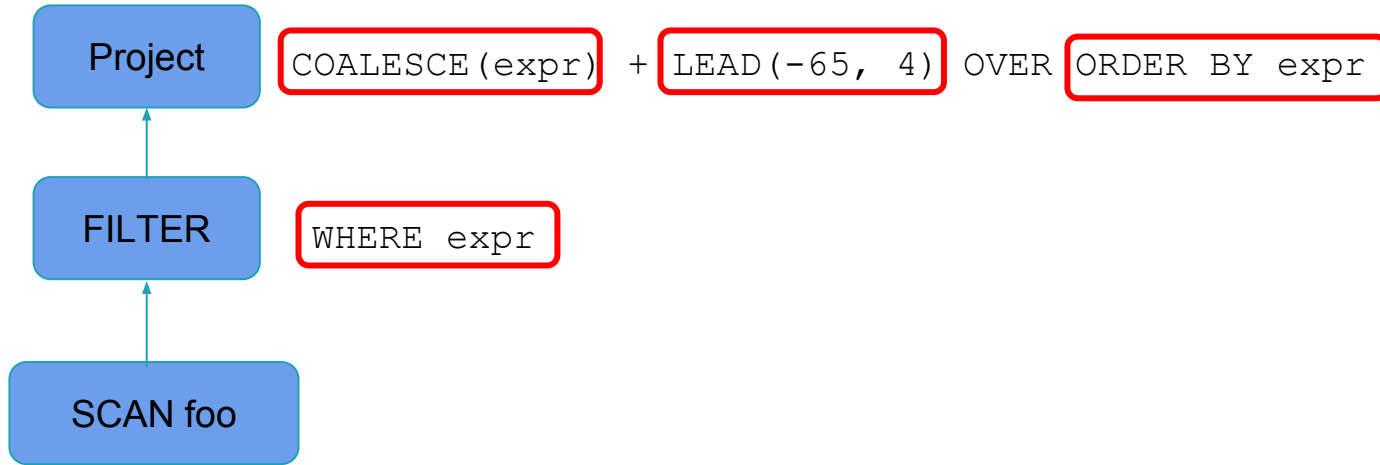
[1.0, 696, -871.81, <<-64.98>>, -349] SPARK row

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

databricks

# 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



Events → Classify

Failure → Correctness

Regression → **Performance** → Re-test → Root-cause → Alert

# 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
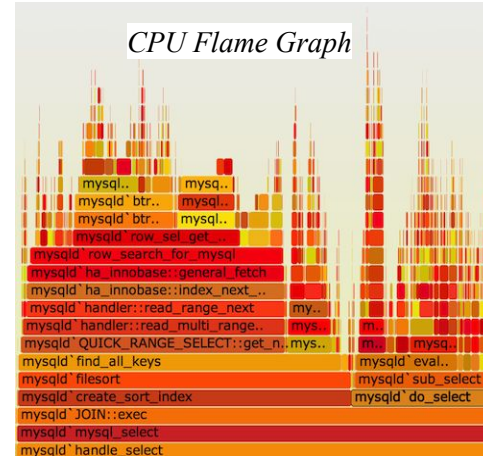
📄 README.md

## 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
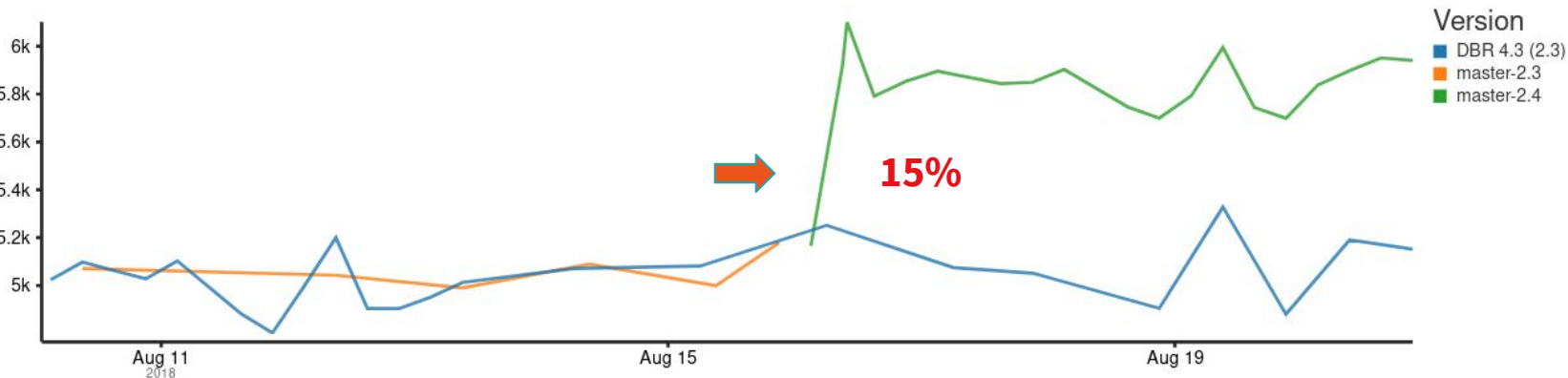


*CPU Flame Graph*

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

daysBack : 30    daysTargetBack : 30    targetVersion : master-2.4

## Time series of runs by type (10 days)



**15%**

**Version**
- DBR 4.3 (2.3)
- master-2.3
- master-2.4

## Top 5 regressing queries to previous version



queryName
- q78
- q93
- q50
- q64
- q84

## Per run and type geomean

# Per query drill-down: 67

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



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

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

databricks

# Q67 executor profile for Spark 2.4-master



LZ4_c..
[unkn..
net/j..
net/j..
net/j..
net/jpou..
net/jpoun..
java/io/B..
java/io/Bu..
java/io/Da..
org/apache..
org/apache/s..
org/apache/spa..

or..
org/apache/spa..
org/apache/spark/sq..
org/apache/spark/sql/exe..
org/apache/spark/sql/exe..
scala/collection/Iterator$$anon$11.hasNext

or..
org/..
org/..
org/..
org/apache/spark/sql/catalyst/expressions/GeneratedClass$GeneratedIteratorForCodeg..
org/apache/spark/sql/catalyst/expressions/GeneratedClass$GeneratedIteratorForCodegenStage7.agg_do..
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

or..
o..
org/..
org/apach..
org/apach..
o.. org/apache/spark/sql/cata..
org/apache/spark/sql/catalyst/expressions/GeneratedClass$GeneratedIteratorForC..
org/apache/spark/sql/c..

org/apach..
org/apach..
org/apache/sp..

or..
org..
org/apache/s..
org/apache/s..

org/..
org/apa..
org/apa..
org/apa..
org/apac..

or..
org/..
org/a..
org/a..
org/a..
org/ap..
scala/c..

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

**Red** slower  **White** new

**Look for hints:**

- Mem mgmt
- Hashing
- unsafe



Diff-After

Reset Zoom                                                                 Search

Blue = speedup, Red == slowdown

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

```scala
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

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

2.4-master — Sep 15, 2018 4,825.72

Version
- 2.3-master
- 2.4-master

15%

5%

< 0%

DBR and Spark testing and performance are a **continuous effort**

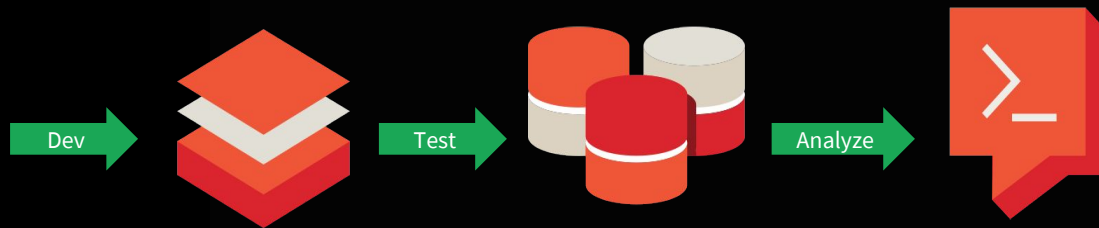- **Over a month** effort to bring performance to improving

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