STITCH FIX

# Spark SQL Performance Troubleshooting and Tuning

Sky Yin

Jun 2, 2016

## Spark is a distributed computing platform

- Distributed =/= Scalable =/= Easy to scale
  - Redshift vs Spark
- Not for data storage (except caching)
- 4 horse chariot
  - 4 language support: Java, Scala, Python, R
  - 4 use cases:
    - SQL
    - Machine learning
    - Graph analysis
    - Stream processing

**Two ways to run a Spark SQL query**

- Flotilla
  - PySpark notebook (not ideal, better to use Presto)
- ETL command
  - Demo 1: https://github.com/stitchfix/quant_workshop/tree/master/sparksql
  - How does it run? ETL => Genie => Spark cluster

## Quick Start

## Resources

- Our [Redshift -> Spark migration guide](#) (thanks Doug, Daragh, Matt, Sky and etc.)
- Spark SQL function reference
    - [`pyspark.sql.functions module`](#) section in PySpark doc
- [Stack Overflow #apache-spark-sql](#)
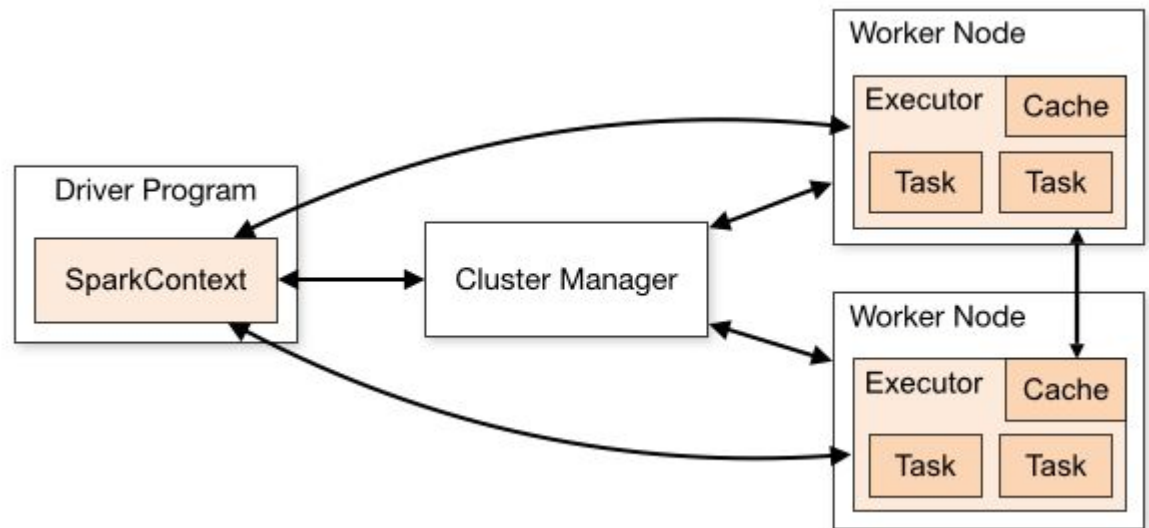- Slack channel #spark-users

# Understanding Spark SQL Execution

## Fundamental building block: RDD

- Resilient Distributed Dataset
- RDD partition =/= Hive metastore partition
- Everything is RDD
    - SQL query

        => a graph of RDD transformation and action

        => Output

# Understanding Spark SQL Execution

## Key question: logical plan => physical plan?

- Demo 2
- Spark execution model
- Physical plan
  - Jobs => Stages => Tasks

# Monitoring Spark SQL Performance

## Spark console
- `go/spark-console` (thanks Tarek)

# Monitoring Spark SQL Performance

## Spark tracking UI

# Monitoring Spark SQL Performance

# Monitoring Spark SQL Performance

# Monitoring Spark SQL Performance

## Details for Stage 0 (Attempt 0)

**Total Time Across All Tasks:** 4.7 min
**Input Size / Records:** 145.6 MB / 11348443
**Shuffle Write:** 155.8 MB / 11348443

▸ DAG Visualization
▸ Show Additional Metrics
▸ Event Timeline

### Summary Metrics for 200 Completed Tasks

| Metric | Min | 25th percentile | Median | 75th percentile | Max |
|---|---|---|---|---|---|
| Duration | | 0.6 s | 0.7 s | 0.9 s | 42 s |
| Scheduler Delay | 4 ms | 7 ms | 9 ms | 12 ms | 0.1 s |
| Task Deserialization Time | 2 ms | 4 ms | 5 ms | 10 ms | 4 s |
| GC Time | | 0 ms | 0 ms | 0 ms | 0.6 s |
| Result Serialization Time | 0 ms | 0 ms | 0 ms | 0 ms | 1 ms |
| Getting Result Time | 0 ms | 0 ms | 0 ms | 0 ms | 0 ms |
| Peak Execution Memory | 0.0 B | 0.0 B | 0.0 B | 0.0 B | 0.0 B |
| Input Size / Records | 0.0 B / 56022 | 758.5 KB / 56580 | 760.4 KB / 56735 | 762.5 KB / 56891 | 769.7 KB / 57550 |
| Shuffle Write Size / Records | 788.9 KB / 56022 | 795.3 KB / 56580 | 797.6 KB / 56735 | 799.8 KB / 56891 | 809.5 KB / 57550 |

## Parallelism

- Number of executor `--spark_num_executors`
- Number of cores/executor `--spark_executor_cores`
- Number of partition
  - Spark will run one task for each RDD partition
  - In PySpark `rdd.getNumPartitions()`
- Repartition

# Tuning Spark SQL Performance

## Memory

- driver vs worker
  - don't call `collect()` on a large query result
- executor memory
  - executor loss
  - `java.lang.OutOfMemoryError: GC overhead limit exceeded`
- GC
- Caching
  - space/speed tradeoff

# Tuning Spark SQL Performance

## Join performance

- Try to avoid `shuffle`
  - `ShuffledHashJoin` vs `BroadcastHashJoin`
  - `spark.sql.autoBroadcastingJoinThreshold`
- `Filter` before `join`
- `Demo 3`
  - [SortMergeJoin](#)
- Detect uneven `shuffle`
  - Tasks taking much longer to run
  - Tasks with much higher input or shuffle output
- Think about the actual distribution of data
  - Ex. group by gender?