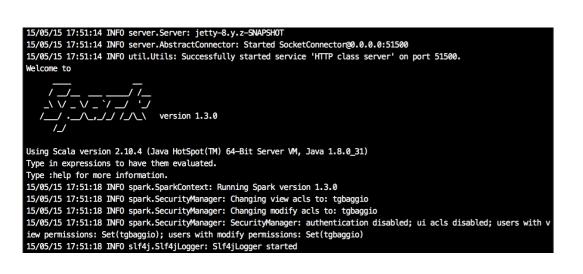# Monitoring and tuning Apache Spark
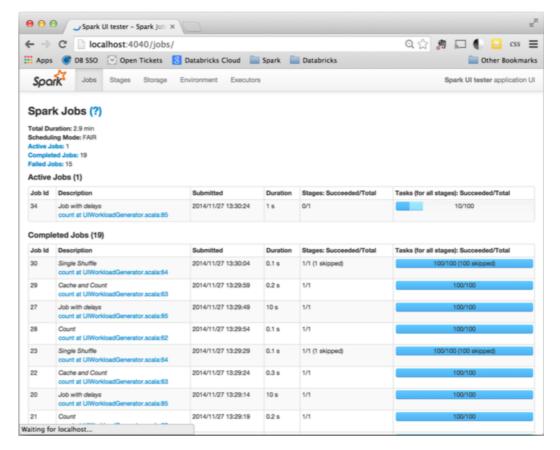
TANG Gen

# Outline

- Components of execution

  - Jobs, stages and tasks

  - Evaluation of DAGs

  - Spark UI

- Spark UI

- Spark logs

- Key performance considerations

# Components of execution
## /Jobs, stages and tasks

- How does a user program get translated into units of physical execution: jobs, stages and tasks ?

# Components of execution
## /Jobs, stages and tasks

- Consider a very simple example
  - Read a log file
  - Split into words and remove empty lines
  - Extract the log level and do a count
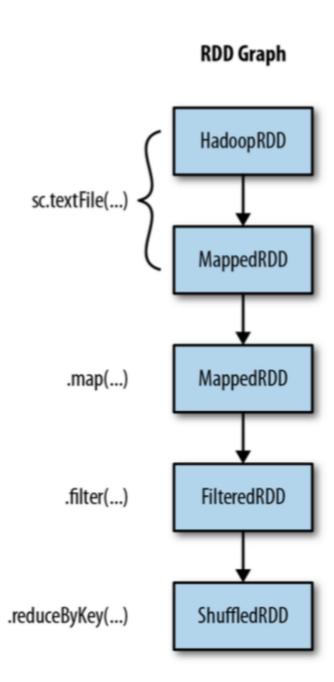
```scala
val input = sc.textFile("log.txt")

val tokenized = input
  .filter(line => line.size > 0)
  .map(line => line.split(" "))

val counts = tokenized
  .map(words => (words(0), 1)
  .reduceByKey({case (a, b) => a + b})
```
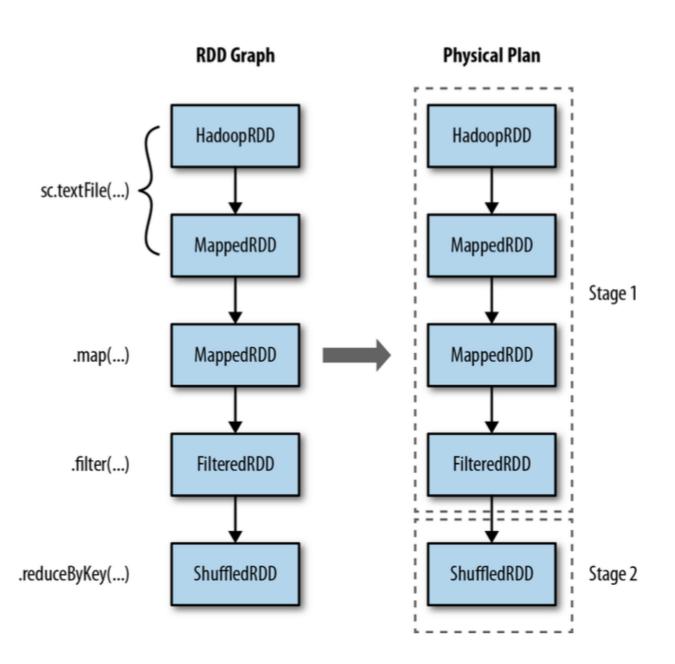
# Components of execution

## /Jobs, stages and tasks

- If we execute the above commands in Spark shell, the program has not performed any actions

- Transformations build up a **DAG**, but don't "do anything

- DAG is directed acyclic graph which will remember the lineage of RDD, including parents, dependencies, etc. It provides fault-tolerance of Spark

- Spark's scheduler creates a physical execution plan according to DAG when an action runs on RDD

**RDD Graph**

HadoopRDD

sc.textFile(...)

MappedRDD

.map(...)     MappedRDD

.filter(...)     FilteredRDD

.reduceByKey(...)     ShuffledRDD

# Components of execution

## /Jobs, stages and tasks

- If we add an action to RDD counts, for example, counts.count()

- Spark's scheduler starts at the final RDD being computed and works backward to find what it much compute

- The scheduler outputs a computation stages for DAG. Each stage has tasks for each partition in RDD. Those stages are executed in reverse order

**RDD Graph**

**Physical Plan**

HadoopRDD

sc.textFile(...)

MappedRDD

.map(...) → MappedRDD

.filter(...) FilteredRDD

.reduceByKey(...) ShuffledRDD

HadoopRDD

MappedRDD

Stage 1

MappedRDD

FilteredRDD

ShuffledRDD Stage 2

# Components of execution
## /Evaluation of DAGs

- Before, we say that "actions" force the evaluation of RDD which is True. But Let's forget it for a moment

- However, in fact, DAGs are materialized through a method sc.runJob in Spark Core

RDD to compute

```scala
def runJob[T, U](
    rdd: RDD[T],
    partitions: Seq[Int],
    func: (Iterator[T] => U)
):Array[U]
```

Which partitions

Function to produce results

# Components of execution

/Evaluation of DAGs

- RunJob needs to compute RDD parents, parents' parents and etc, util to an RDD with no dependencies

**runJob(counts)**
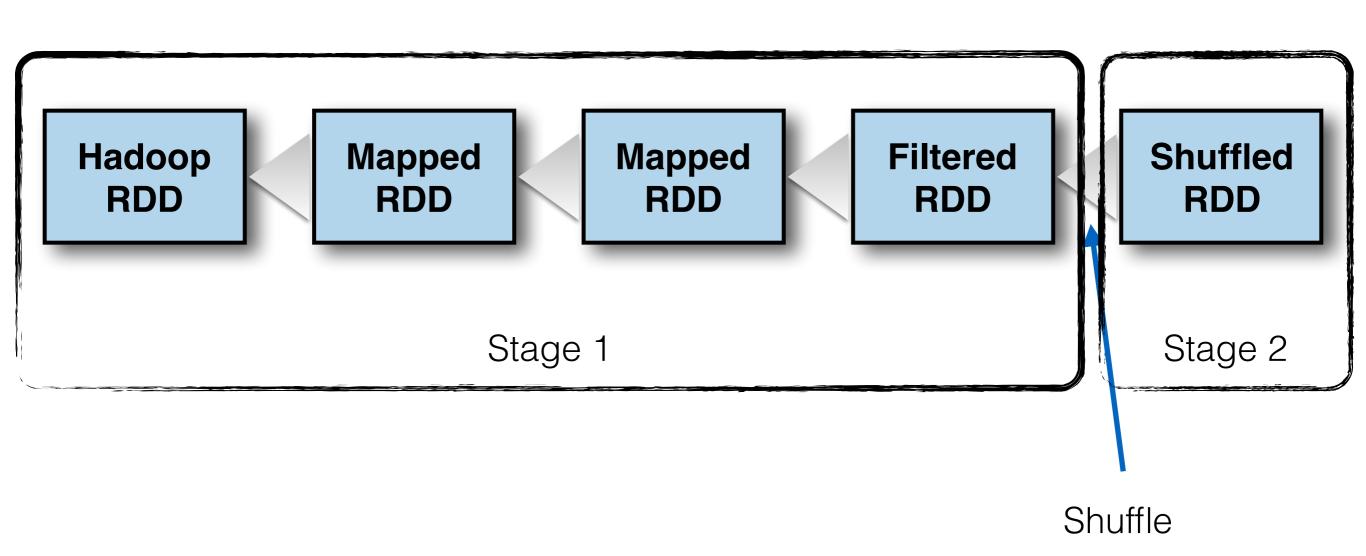
| **Hadoop RDD** | | **Mapped RDD** | | **Mapped RDD** | | **Filtered RDD** | | **Shuffled RDD** |

Input                                    tokenized                                    counts

# Components of execution
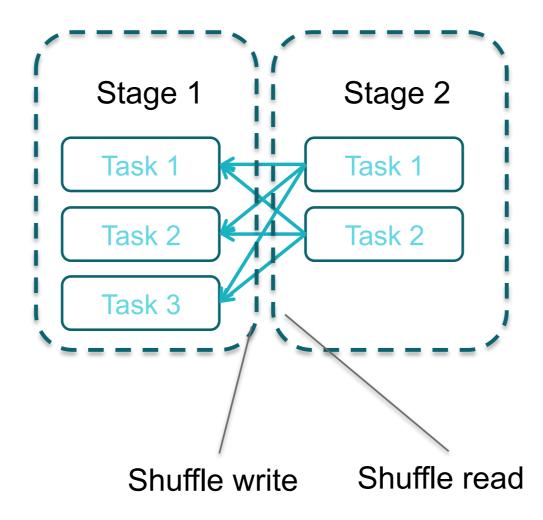## /Physical optimization

- Certain types of transformation can be **pipelined**

- If dependent RDD's have already been cached (or persisted in a shuffle), the graph can be **truncated**

- Once pipelining and truncation occur, Spark produces a set of **stages**, each stage is composed of **tasks**

# Components of execution

## /Physical optimization

# Components of execution
## /Stage graph

- In stage 1, each task will

  1. Read hadoop file

  2. Perform maps and filters

  3. Write partial sums

**Stage 1**

Task 1

Task 2

Task 3

**Stage 2**

Task 1

Task 2

Shuffle write    Shuffle read

- In stage 2, each task will

  1. Read partial sums

  2. Invoke user function passed to runJob

# Components of execution

/Count() action

```
class RDD {
  def count(): Long = {
    results = sc.runJob(
      this,
      0 until partitions.size
      it => it.size()
    )
  return results.sum
  }
}
```

RDD = self

Partitions = all partitions

Function = size of all partition

# Components of execution
## /Spark UI

localhost:4040/jobs/job/?id=0

Apps | DB SSO | Open Tickets | Databricks Cloud | Spark | Databricks | Databricks Root Fold | Shard Listings | Other Bookmarks

**Spark** 1.3.0-SNAPSHOT | **Jobs** | Stages | Storage | Environment | Executors | **Spark shell** application

### Details for Job 0

**Status:** SUCCEEDED
**Completed Stages:** 2

**Completed Stages (2)**

| Stage Id | Description | | Submitted | Duration | Tasks: Succeeded/Total | Input | Output | Shuffle Read | Shuffle Write |
|---|---|---|---|---|---|---|---|---|---|
| 1 | count at <console>:28 | +details | 2015/02/19 21:35:56 | 38 ms | 2/2 | | | | |
| 0 | map at <console>:25 | +details | 2015/02/19 21:35:56 | 0.1 s | 2/2 | 72.0 B | | | 354.0 B |

Named after action calling runJob

Named after last RDD in pipeline

13

# Components of execution
## /Summary

- Definition:

    - **Jobs**: Work required to compute RDD in runJob.

    - **Stages:** A wave of work within a job, corresponding to one or more pipelined RDD's.

    - **Tasks:** A unit of work within a stage, corresponding to one RDD partition.

    - **Shuffle:** The transfer of data between stages.

- Phases occur during Spark execution:

    - User code defines a DAG of RDDs

    - Actions force translation of the DAG to an execution plan

    - Tasks are scheduled and executed on a cluster

# Spark UI
## /Demo time

- If Spark runs on standalone or Mesos cluster, Spark's built-in web UI is available on the machine where driver is running at port 4040 by default

- If Spark runs on YARN cluster mode，the driver runs inside the cluster, the UI is accessed through the YARN ResourceManager

# Spark logs
## /Driver and executor logs

- The location of Spark's log files depends on the deployment mode:
  - Standalone, application logs are directly displayed in the stand-alone master's web UI. They are stored by default in the *work/* directory of the Spark distribution on each worker.
  - In Mesos, logs are stored in the *work/* directory of a Mesos slave, and accessible from Mesos master UI
  - In YARN mode, when the application finishes, use the following command to produce a report containing logs

    *yarn logs -applicationID <app ID>*

  - In YARN mode, when the application is running, we can access the logs of certain containers via ResourceManager UI

# Spark logs
## /Logging system

- Spark's logging subsystem is based on *log4j,* a widely used Java library

- An example log4j configuration file is bundled with Spark at *conf/ log4j.properties.template*. To customize Spark's logging,

    1. Copy the example to a file called *log4j.properties.*

    2. Modify behavior such as the root logging level (the threshold level for logging output).

    3. Add the *log4j.properties* file using the *--files* flag of *spark-submit*

# Key performance considerations

/Level of Parallelism

- RDD's partitions size decides the level of parallelism during execution, which can be modified by *coalesce(Num)* or *repartition(Num)*

- Demo time

  - Too little parallelism: idle cores

  - Too much parallelism: too much shuffle

# Key performance considerations
/Serialization format

- When Spark is transferring data over the network or spilling data to disk, it needs to serialize objects into a binary format

- By default, Spark will use Java's built-in serializer.

- However, Spark also support the use of Kyro, a third-party serialization library (Kyro is used in mllib library)

- To use Kyro serializer

```
val conf = SparkConf()
conf.set("spark.serializer", "org.apache.spark.serializer.KyroSerializer")
conf.set("spark.kyro.registrationRequired", "true")
conf.registerKyroClasses(Array(classOf[Myclass], classOf[Myotherclass]))
```

# Key performance considerations

## /Memory management

- In general, memory is used for the following three ways:

  - RDD storage

  - Shuffle and aggregation buffers

  - User code

- The more detailed presentation of memory management will be in tommorrow