

# Data Engineering with Apache Spark

---

Markus Dale, [medale@asymmetrik.com](mailto:medale@asymmetrik.com)

May 2019

- Slides: <https://github.com/medale/>
- Scala Spark Code Examples: <https://github.com/medale/>

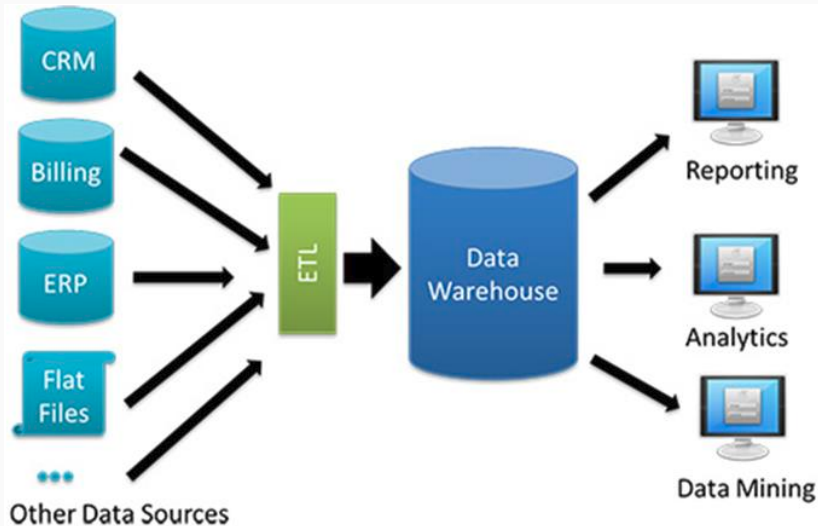




Figure 1: Laptop

## Data engineering for larger dataset (Vertical Scaling)

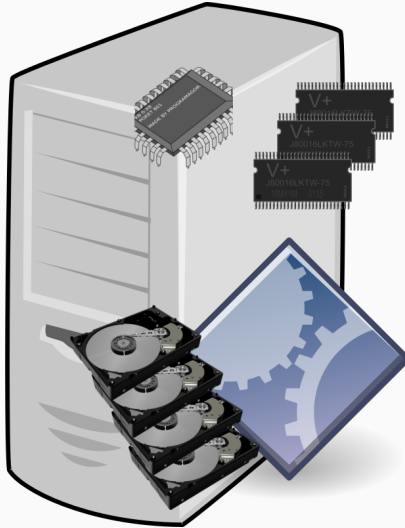


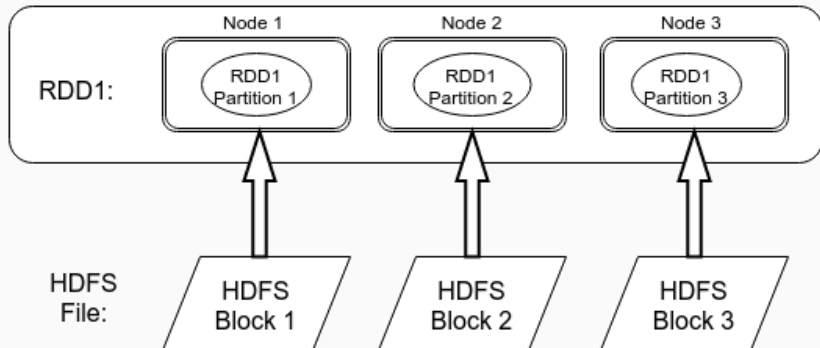
Figure 2: Beefed-up Server

## Data engineering for large datasets (Horizontal Scaling)



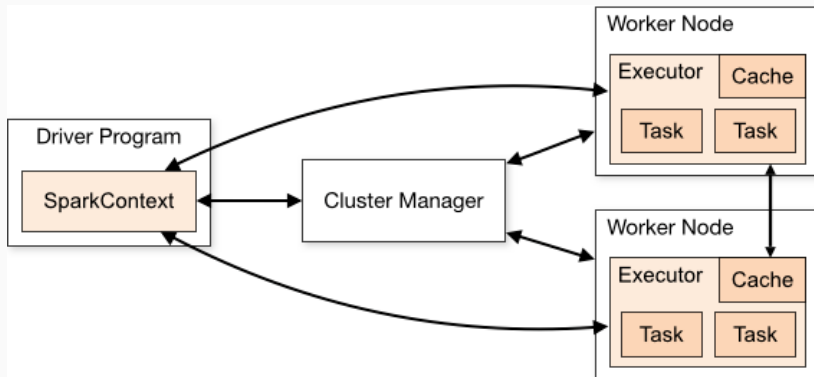


## Resilient Distributed Datasets (RDDs)





# Anatomy of a Spark Application



Source: Apache Spark website

# Hello, Spark World!

```
import org.apache.spark.sql.SparkSession

object HelloSparkWorld {

  def process(spark: SparkSession): (Long,Long) = {
    val records = spark.read.json( path = "file:///datasets/github/data")
    records.cache()
    val totalEventCount = records.count()

    val prs = records.where(records("type") === "PullRequestEvent")
    val pullRequestEventCount = prs.count()

    records.unpersist()
    (totalEventCount, pullRequestEventCount)
  }

  def main(args: Array[String]): Unit = {
    val spark = SparkSession.builder().
      appName( name = "HelloSparkWorld").
      getOrCreate()
    process(spark)
  }
}
```

## Starting Spark Standalone Cluster Manager

# Start on master

```
$SPARK_HOME/sbin/start-master.sh --host 192.168.1.230
```

# Start one or more workers

```
$SPARK_HOME/sbin/start-slave.sh spark://192.168.1.230:7077
```

# Spark Standalone Cluster Manager UI - idle



## Spark Master at spark://192.168.1.230:7077

URL: spark://192.168.1.230:7077

Alive Workers: 1

Cores in use: 8 Total, 0 Used

Memory in use: 30.4 GB Total, 0.0 B Used

Applications: 0 [Running](#), 0 [Completed](#)

Drivers: 0 Running, 0 Completed

Status: ALIVE

### ▼ Workers (1)

Worker Id	Address	State	Cores	Memory
<a href="#">worker-20190430220608-192.168.1.230-37667</a>	192.168.1.230:37667	ALIVE	8 (0 Used)	30.4 GB (0.0 B Used)

### ▼ Running Applications (0)

Application ID	Name	Cores	Memory per Executor	Submitted Time	User	State	Duration
----------------	------	-------	---------------------	----------------	------	-------	----------

### ▼ Completed Applications (0)

Application ID	Name	Cores	Memory per Executor	Submitted Time	User	State	Duration
----------------	------	-------	---------------------	----------------	------	-------	----------

```
spark-shell --master spark://192.168.1.230:7077 \  
  --driver-memory 1g \  
  --executor-memory 2g \  
  --total-executor-cores 4 \  
  --executor-cores 2 \  
  --jars /tmp/dataset-0.9.0-SNAPSHOT-fat.jar
```

# Spark Standalone Cluster Manager - 1 running application



## Spark Master at spark://192.168.1.230:7077

URL: spark://192.168.1.230:7077

Alive Workers: 1

Cores in use: 8 Total, 4 Used

Memory in use: 30.4 GB Total, 4.0 GB Used

Applications: 1 [Running](#), 0 [Completed](#)

Drivers: 0 Running, 0 Completed

Status: ALIVE

### Workers (1)

Worker Id	Address	State	Cores	Memory
<a href="#">worker-20190430220608-192.168.1.230-37667</a>	192.168.1.230:37667	ALIVE	8 (4 Used)	30.4 GB (4.0 GB Used)

### Running Applications (1)

Application ID	Name	Cores	Memory per Executor	Submitted Time	User	State	Duration
<a href="#">app-20190430221543-0000</a> (kill)	<a href="#">Spark shell</a>	4	2.0 GB	2019/04/30 22:15:43	medale	RUNNING	14 min

### Completed Applications (0)

Application ID	Name	Cores	Memory per Executor	Submitted Time	User	State	Duration
----------------	------	-------	---------------------	----------------	------	-------	----------

## RDDs - Not deprecated!

```
object RddProcessor {  
  
  val DefaultEventInputUrl = "file:///datasets/github/data"  
  
  def process(sc: SparkContext, inputUrl: String): (Long, Long) = {  
    val records: RDD[String] = sc.textFile(inputUrl)  
    println(s"We have a total of ${records.partitions.size} partitions.")  
    val total = records.count()  
    val prs = records.filter(r => r.contains("PullRequestEvent"))  
    val totalPrs = prs.count()  
    (total, totalPrs)  
  }  
  
  def main(args: Array[String]): Unit = {  
    val spark = SparkSession.builder().  
      appName(name = "RddProcessor").  
      getOrCreate()  
  
    val inputUrl = if (args.size > 0) {  
      args(0)  
    } else {  
      DefaultEventInputUrl  
    }  
    process(spark.sparkContext, inputUrl)  
    spark.stop()  
  }  
}
```

## And now for something completely different: Colon Cancer



- Screening saves lives!
  - Colonoscopy - talk to your doc
  - Dave Barry: A journey into my colon — and yours
- Colorectal Cancer Alliance



# Questions?



- medale@asymmetrik.com
- Infrequent blog/past presentations <http://uebercomputing.com/>