

# Apache Spark Through Email

---

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

July 2024

- Slides: <https://github.com/medale/spark-mail/blob/master/presentation/ApacheSparkThroughEmail.pdf>
- Spark Code Examples: <https://github.com/medale/spark-mail/>
  - README.md describes how to get and parse Enron email dataset

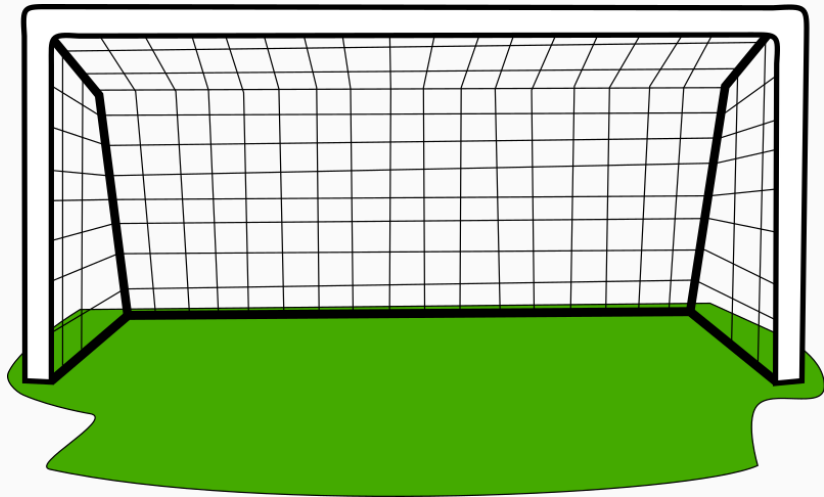


Figure 1: Intro to Apache Spark



Figure 2: Laptop

## Data Science for Larger Dataset (Vertical Scaling)

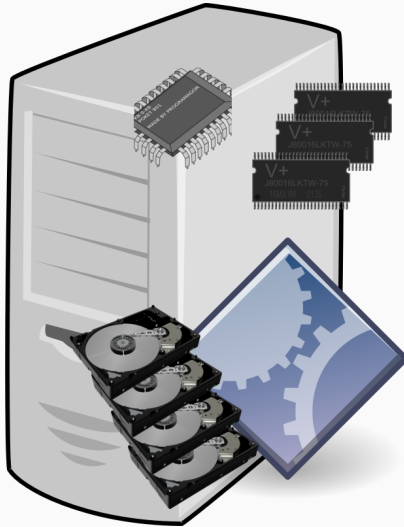


Figure 3: Beefed-up Server

## Data Science for Large Datasets (Horizontal Scaling)





Figure 5: HDFS, MapReduce

# Hadoop Ecosystem

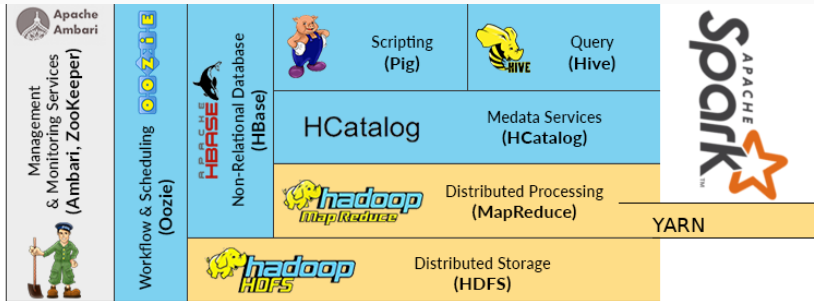


Figure 6: Some Frameworks Around Hadoop



- Local
  - Download from <https://spark.apache.org>, untar, add to PATH
  - SDKMAN - `curl -s "https://get.sdkman.io" | bash`
    - `sdk install spark`
  - `spark-shell` or `pyspark`
- Standalone cluster, Hadoop YARN
  - Need shared file system or common datastore (e.g. AWS S3)
- Cloud-based managed:
  - AWS EMR
  - GCP Dataproc
  - Databricks on Azure, GCP or AWS

- Scala
- Java
- Python
- R
- => Project Tungsten code generation

# Apache Spark Components

Structured  
Streaming

Advanced  
Analytics

Libraries &  
Ecosystem

Structured APIs

Datasets

DataFrames

SQL

Low-level APIs

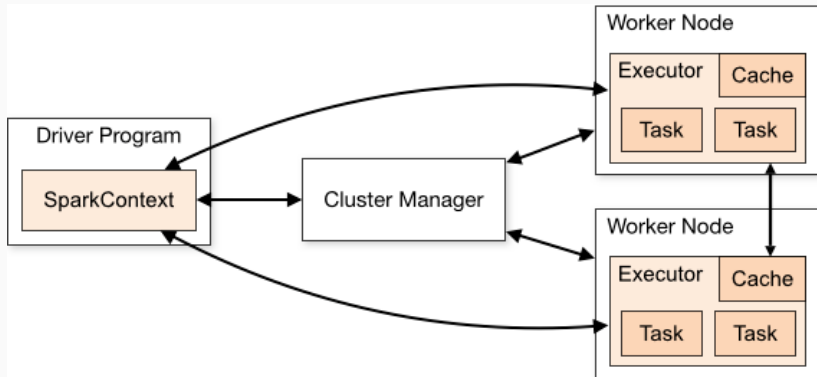
RDDs

Distributed Variables

Source: Spark: The Definitive Guide

- Jupyter Notebook with Apache Toree
- See Notebook [../notebooks/html/SparkThroughEmail1.html](http://../notebooks/html/SparkThroughEmail1.html)

# Cluster Manager, Driver, Executors, Tasks



Source: Apache Spark website

## SparkSession: Entry to cluster

- spark: spark.sql.SparkSession

```
//SparkSession provided by notebook or shell as spark
val homeDir = sys.props("user.home")
val records = spark.read.
    parquet(s"$homeDir/datasets/enron/enron-small.parquet")

//In regular code for spark-submit
//com.uebercomputing.spark.dataset.TopNEmailMessageSenders
val spark = SparkSession.builder().
    appName("TopNEmailMessageSenders").
    master("local[2]").getOrCreate()
```

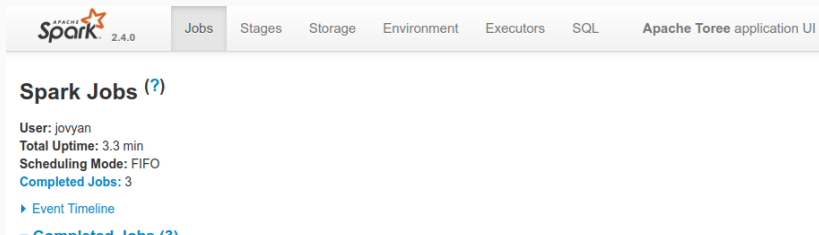
- `spark.read`: `spark.sql.DataFrameReader`
  - jdbc
  - json
  - parquet
  - text...
  - Also: <https://spark-packages.org> - Avro, Redshift, MongoDB...

# Transformations vs. Actions

- Transformation: returns a new RDD (nothing gets executed)
  - `read`, `cache`, `filter`...
- Actions: trigger execution, catalyst query optimizer, Tungsten code generation
  - `count`, `write`, ...

## Scaling Behind the Scenes

---



The screenshot shows the Apache Spark Jobs page within the Apache Torte application UI. The top navigation bar includes tabs for Jobs, Stages, Storage, Environment, Executors, SQL, and the Apache Torte application UI. The main content area displays the following information:

- Spark Jobs (?)**
- User:** jovyan
- Total Uptime:** 3.3 min
- Scheduling Mode:** FIFO
- Completed Jobs:** 3
- [▶ Event Timeline](#)
- [Completed Jobs \(?\)](#)



# Stages: Pipeline work per stage - shuffle

## ▼ DAG Visualization

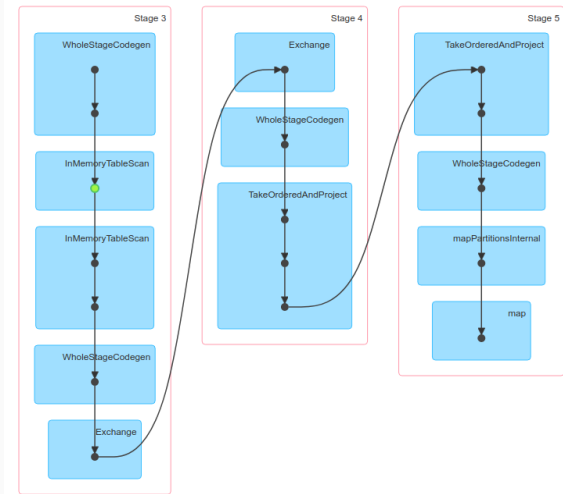


Figure 8: Stages


- See Notebook [../notebooks/html/ApacheSparkThroughEmail2.html](http://notebooks/html/ApacheSparkThroughEmail2.html)

# Spark APIs - DataFrameReader, Dataset, Column, functions

▶	<code>def <a href="#">parquet</a>(paths: String*): <a href="#">DataFrame</a></code> Loads a Parquet file, returning the result as a <a href="#">DataFrame</a> .
▶	<code>def <a href="#">parquet</a>(path: String): <a href="#">DataFrame</a></code> Loads a Parquet file, returning the result as a <a href="#">DataFrame</a> .
▶	<code>def <a href="#">schema</a>(schemaString: String): <a href="#">DataFrameReader</a></code> Specifies the schema by using the input DDL-formatted string.
▶	<code>def <a href="#">schema</a>(schema: <a href="#">StructType</a>): <a href="#">DataFrameReader</a></code> Specifies the input schema.
▶	<code>def <a href="#">table</a>(tableName: String): <a href="#">DataFrame</a></code> Returns the specified table as a <a href="#">DataFrame</a> .
▶	<code>def <a href="#">text</a>(paths: String*): <a href="#">DataFrame</a></code> Loads text files and returns a <a href="#">DataFrame</a> whose schema starts with a string column named "value".
▶	<code>def <a href="#">text</a>(path: String): <a href="#">DataFrame</a></code> Loads text files and returns a <a href="#">DataFrame</a> whose schema starts with a string column named "value".
▶	<code>def <a href="#">textFile</a>(paths: String*): <a href="#">Dataset[String]</a></code> Loads text files and returns a <a href="#">Dataset</a> of <a href="#">String</a> .
▶	<code>def <a href="#">textFile</a>(path: String): <a href="#">Dataset[String]</a></code> Loads text files and returns a <a href="#">Dataset</a> of <a href="#">String</a> .

Expression operators	
▶	<code>def <a href="#">%(other: Any): <a href="#">Column</a></a></code> Modulo (a.k.a.
▶	<code>def <a href="#">%(other: Any): <a href="#">Column</a></a></code> Boolean AND.
▶	<code>def <a href="#">*(other: Any): <a href="#">Column</a></a></code> Multiplication of this expression and another expression.
▶	<code>def <a href="#">+(other: Any): <a href="#">Column</a></a></code> Sum of this expression and another expression.
▶	<code>def <a href="#">-(other: Any): <a href="#">Column</a></a></code> Subtraction.
▶	<code>def <a href="#">/(other: Any): <a href="#">Column</a></a></code> Division this expression by another expression.
▶	<code>def <a href="#">&lt;(other: Any): <a href="#">Column</a></a></code> Less than.
▶	<code>def <a href="#">&lt;=(other: Any): <a href="#">Column</a></a></code> Less than or equal to.
▶	<code>def <a href="#">==(other: Any): <a href="#">Column</a></a></code> Equality test that is safe for null values.
▶	<code>def <a href="#">!=(other: Any): <a href="#">Column</a></a></code> Inequality test.

Date time functions	
▶	<code>def <a href="#">add_months</a>(startDate: <a href="#">Column</a>, numMonths: Int): <a href="#">Column</a></code> Returns the date that is numMonths after startDate.
▶	<code>def <a href="#">current_date</a>(): <a href="#">Column</a></code> Returns the current date as a date column.
▶	<code>def <a href="#">current_timestamp</a>(): <a href="#">Column</a></code> Returns the current timestamp as a timestamp column.
▶	<code>def <a href="#">date_add</a>(start: <a href="#">Column</a>, days: Int): <a href="#">Column</a></code> Returns the date that is days days after start
▶	<code>def <a href="#">date_format</a>(dateExpr: <a href="#">Column</a>, format: String): <a href="#">Column</a></code> Converts a date/timestamp to a value of string in the format specified by the second argument.
▶	<code>def <a href="#">date_sub</a>(start: <a href="#">Column</a>, days: Int): <a href="#">Column</a></code> Returns the date that is days days before start
▶	<code>def <a href="#">date_trunc</a>(format: String, timestamp: <a href="#">Column</a>): <a href="#">Column</a></code> Returns timestamp truncated to the unit specified by the format.
▶	<code>def <a href="#">datediff</a>(end: <a href="#">Column</a>, start: <a href="#">Column</a>): <a href="#">Column</a></code> Returns the number of days from start to end
▶	<code>def <a href="#">dayofmonth</a>(e: <a href="#">Column</a>): <a href="#">Column</a></code> Extracts the day of the month as an integer from a given date/timestamp/string.
▶	<code>def <a href="#">dayofweek</a>(e: <a href="#">Column</a>): <a href="#">Column</a></code> Extracts the day of the week as an integer from a given date/timestamp/string.

 Dataset

`class Dataset[T] extends Serializable`

A [Dataset](#) is a strongly typed collection of domain-specific objects that can be transformed in parallel using functional or relational operations. Each operation available on [Datasets](#) is divided into transformations and actions. Transformations are the ones that produce new [Datasets](#), and actions filter, select, and aggregate (groupBy). Example actions count, show, or writing data out to the filesystem.

[Datasets](#) are "lazy" (i.e. computations are only triggered when an action is invoked). Internally, a [Dataset](#) represents a logical plan that describes the optimized logical plan and generates a physical plan for efficient execution in a parallel and distributed manner. To explore the logical plan on which Spark to generate code at runtime to serialize the [Parquet](#) object into a binary structure. This binary structure often has much lower memory overhead than the internal binary representation for data, use the [Schema](#) function.

There are typically two ways to create a [Dataset](#). The most common way is by pointing Spark to some files on storage systems, using the `read` function.

```
val people = spark.read.parquet("...").as[Person] // Scala
Dataset<Person> people = spark.read().parquet("...").as(Encoders.bean(Person.class)); // Java
```

[Datasets](#) can also be created through transformations available on existing [Datasets](#). For example, the following creates a new [Dataset](#) by applying transformations to an existing [Dataset](#).

```
val names = people.map(_._name) // in Scala; names is a Dataset[String]
Dataset<String> names = people.map((Person p) => p.name, Encoders.STRING());
```

[Dataset](#) operations can also be untyped, through various domain-specific language (DSL) functions defined in [Dataset](#) (this class), [Column](#), and [SQL](#) in R or Python.

To select a column from the [Dataset](#), use `apply` method in Scala and `col` in Java.

```
val ageCol = people["age"] // in Scala
Column ageCol = people.col("age"); // in Java
```

Note that the [Column](#) type can also be manipulated through its various functions.

```
// The following creates a new column that increases everybody's age by 10.
people["age"] + 10 // in Scala
people.col("age").plus(10); // in Java
```

- Goldilocks - not too many, not too few
- Initial parallelism - number of input “blocks”
- Shuffle - `spark.sql.shuffle.partitions` configuration

- See Notebook [../notebooks/html/ApacheSparkThroughEmail3.html](http://notebooks/html/ApacheSparkThroughEmail3.html)

- See <https://spark.apache.org/pandas-on-spark/>

- Avro - record-oriented data format
- Parquet - column-oriented data format by page
- Arrow - share memory for Python

([https://spark.apache.org/docs/latest/api/python/user\\_guide/sql/arrow\\_pa](https://spark.apache.org/docs/latest/api/python/user_guide/sql/arrow_pa)

- <https://spark.apache.org/>
- <https://spark-packages.org/> - Community 3rd party packages (e.g. data sources)
- <https://sparkbyexamples.com/>



## And now for something completely different: Colon Cancer



- Screening saves lives!
  - Colonoscopy - talk to your doc
- Colorectal Cancer Alliance

# Questions?



- markus.dale@bluehalo.com
- Infrequent blog/past presentations <http://uebercomputing.com/>
- Spark Mail repo <https://github.com/medale/spark-mail/>