# Apache Spark Through Email

Markus Dale, 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

Figure 3: Beefed-up Server

**Figure 5:** HDFS, MapReduce

Figure 6: Some Frameworks Around Hadoop

## Running Spark

- 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`
    - Edit `$SPARK_HOME/conf/spark-defaults.conf` (from template)
        - `spark.driver.memory              8g`
- 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

Structured Streaming

Advanced Analytics

Libraries & Ecosystem

Structured APIs

Datasets          DataFrames          SQL
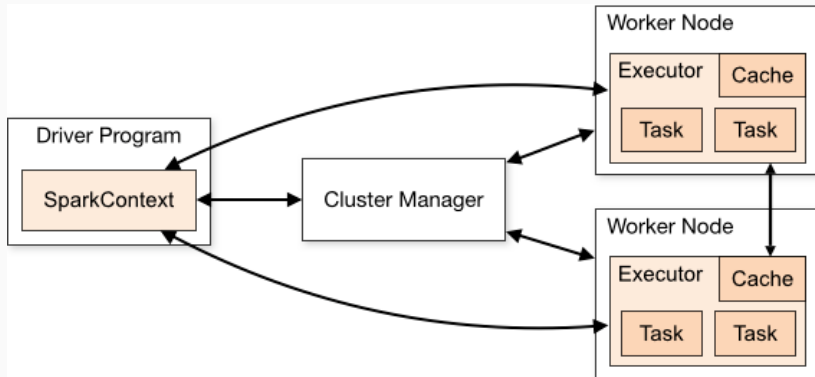
Low-level APIs

RDDs          Distributed Variables

Source: Spark: The Definitive Guide

```
spark-shell --master "local[4]" --driver-memory 8G
```

- Jupyter notebook with Apache Toree Notebook
  ../notebooks/html/ApacheSparkThroughEmail1.html

Source: Apache Spark website

- **spark**: spark.sql.SparkSession

```scala
//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/write`: spark.sql.DataFrameReader/Writer
    - jdbc
    - json
    - parquet
    - text...
    - Also: https://spark-packages.org - Redshift, MongoDB...

## Convert Dataset Format

```scala
import org.apache.spark.sql.functions._
val homeDir = sys.props("user.home")
val records = spark.read.parquet(s"$homeDir/datasets/enron/enron-small.parquet")

// write block-size file(s)
records.write.json(s"$homeDir/datasets/enron/json-parts")
records.repartition(1).write.json(s"$homeDir/datasets/enron/json-single")

// Dataset has mailfields/RE: and mailfields/re: fields
spark.conf.set('spark.sql.caseSensitive', true)
val jsonIn = spark.read.json(s"$homeDir/datasets/enron/json-parts")
```

- Transformation: returns a new RDD (nothing gets executed)
  - `read`, `cache`, `select`, `where`…
- Actions: trigger execution, catalyst query optimizer, Tungsten code generation
  - `count`
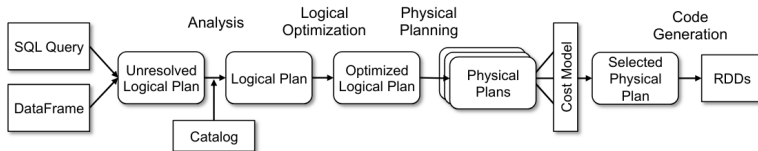  - Bring rows back to driver: `take`, `collect` (watch OOM!)
  - `write`

Figure 3: Phases of query planning in Spark SQL. Rounded rectangles represent Catalyst trees.

**Figure 7:** Jobs and Tasks

Figure 8: Stages

- See Notebook ../notebooks/html/ApacheSparkThroughEmail2.html

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

Operations available on Datasets are divided into transformations and actions. Transformations are the ones that produce new Datasets, and actions are the ones that trigger computation and return results. Example transformations include map, filter, select, and aggregate (groupBy). Example actions count, show, or writing data out to file systems.

Datasets are "lazy", i.e. computations are only triggered when an action is invoked. Internally, a Dataset represents a logical plan that describes the computation required to produce the data. When an action is invoked, Spark's query optimizer optimizes the logical plan and generates a physical plan for efficient execution in a parallel and distributed manner. To explore the logical plan as well as

To efficiently support domain-specific objects, an *Encoder* is required. The encoder maps the domain specific type T to Spark's internal type system. For example, given a class Person with two fields, name (string) and age (int), an encoder is used to tell Spark to generate code at runtime to serialize the Person object into a binary structure. This binary structure often has much lower memory footprint as well as are optimized for efficiency in data processing (e.g. in a columnar format). To understand 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 available on a SparkSession.

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

Datasets can also be created through transformations available on existing Datasets. For example, the following creates a new Dataset by applying a filter on the existing one:

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

Dataset operations can also be untyped, through various domain-specific-language (DSL) functions defined in: Dataset (this class), Column, and functions. These operations are very similar to the operations available in the data frame abstraction 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
```

**Expression operators**

| | def **%**(other: Any): <u>Column</u> |
|---|---|
| ▸ | Modulo (a.k.a. |
| ▸ | def **&&**(other: Any): <u>Column</u> |
| | Boolean AND. |
| ▸ | def ***(other: Any): <u>Column</u> |
| | Multiplication of this expression and another expression. |
| ▸ | def **+**(other: Any): <u>Column</u> |
| | Sum of this expression and another expression. |
| ▸ | def **-**(other: Any): <u>Column</u> |
| | Subtraction. |
| ▸ | def **/**(other: Any): <u>Column</u> |
| | Division this expression by another expression. |
| ▸ | def **<**(other: Any): <u>Column</u> |
| | Less than. |
| ▸ | def **<=**(other: Any): <u>Column</u> |
| | Less than or equal to. |
| ▸ | def **<=>**(other: Any): <u>Column</u> |
| | Equality test that is safe for null values. |
| ▸ | def **!=**(other: Any): <u>Column</u> |
| | Inequality test. |

**Date time functions**

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

- Goldilocks - not too many, not too few
- Initial parallelism - number of input "blocks"
- Splittable file formats (e.g. parquet, avro, bzip2)
    - Not zip, gzip!
- Shuffle - Adaptive Query Execution (dynamic partitioning)

- See Notebook ../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

# Resources

- https://spark.apache.org/
- https://spark-packages.org/ - Community 3rd party packages (e.g. data sources)
- https://sparkbyexamples.com/
- RDD - https://www.usenix.org/system/files/conference/nsdi12/nsdi12-final138.pdf
- Spark SQL - https://people.csail.mit.edu/matei/papers/2015/sigmod_spark_sql.pdf
- Adaptive query execution - https://www.databricks.com/blog/2020/05/29/adaptive-query-execution-speeding-up-spark-sql-at-runtime.html

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

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