# EGT 305: BIG DATA PROCESSING & APPLICATION

BOSS

PRESENTATION TITLE

# WHAT IS APACHE SPARK

- Apache Spark is an <mark>open-source distributed computing system</mark> designed for big data processing and analytics.

- Unified analytics engine for large-scale data processing with speed, ease of use, and versatility.

- Spark supports various programming languages including Scala, Java, Python, and R, making it accessible to a wide range of developers.

- Run queries and machine learning workflows on petabytes of data, which is impossible to do on your local device.

2

Source: https://www.oracle.com/sg/big-data/what-is-big-data/#:~:text=What%20exactly%20is%20big%20data,especially%20from%20new%20data%20sources

# HOW APACHE SPARK WORKS? - FROM RDD TO DATAFRAME

**Resilient Distributed Datasets (RDDs)** - low level data representation in Spark

**DataFrame API was introduced in 2015** – RDD + schema + advance interfaces

✓ Incredibly powerful API:

- **Hide low level distributed operations.** Map-reduce tasks that used to take thousands of lines of code to express could be reduced to dozens.

```
spark.textFile("hdfs://...")
.flatMap(line => line.split(" "))
.map(word => (word, 1))
.reduceByKey(_ + _)
.saveAsTextFile("hdfs://...")
```
**VS**
Written in Java for MapReduce it has around 50 lines of code

- Conceptually equivalent to a table in a relational database or a data frame in R/Python, but with richer **optimizations** under the hood.
- Supporting multiple general-purpose programming languages (Java, Python, Scala).

✓ Ability to scale
  **No need to change code:** from kilobytes of data on a single laptop to petabytes on a large cluster**.**

✓ Support not only batch data, but also **streaming data**
  - To write streaming jobs the same way you write batch jobs.
  - Join streams against historical data, or run ad-hoc queries on stream state.
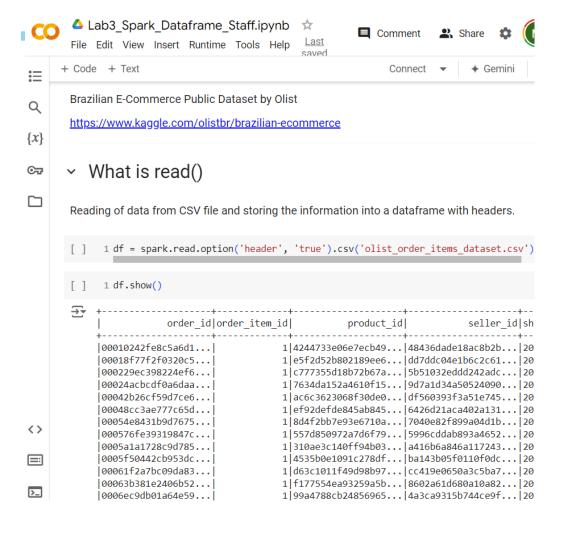  - Build powerful interactive applications, not just analytics.

# LAB 3:

## DATAFRAME

- SPARK DATAFRAME FUNDAMENTALS:
  - read()
  - explain()
  - printSchema()
  - inferSchema()
  - DoubleData Type
  - Other Datattypes
  - Data Transformation
  - Data ETL using PySpark

4

PRESENTATION TITLE

# LAB 3



- Download file from Brightspace
- Transfer file into Google Colab
- Start running the code

5

# THANK YOU

SLEEPYHEAD
SLEPPYMEL@GMAIL.COM