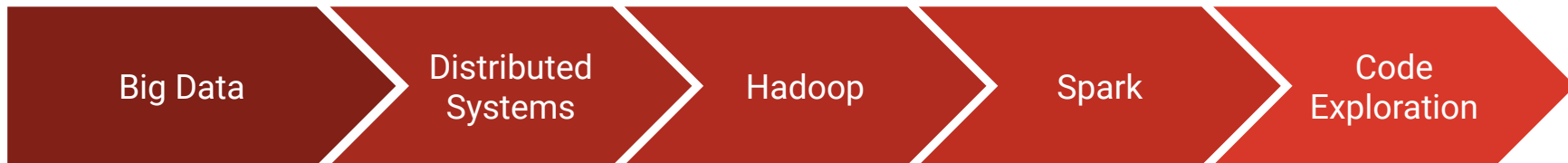


Introduction to Big Data in PySpark

Road Map:



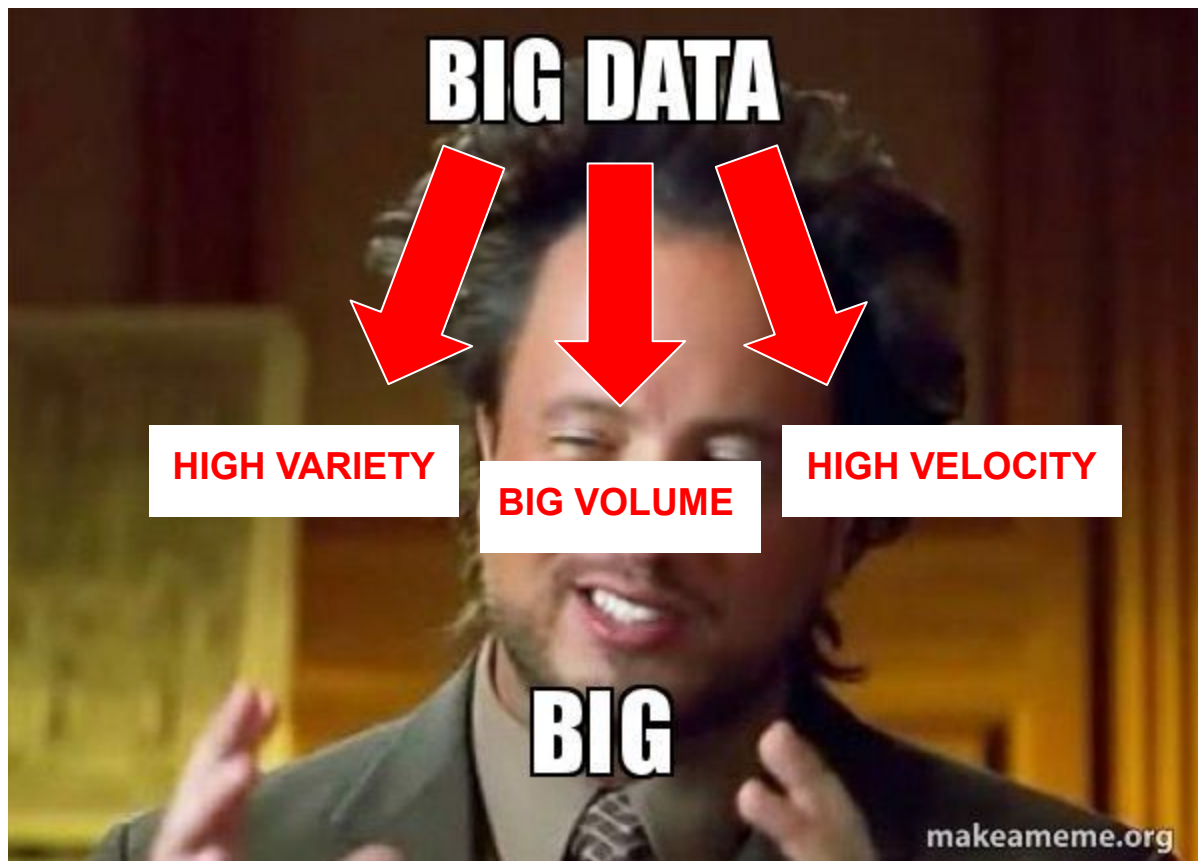
Objectives:

- Introduction to distributed systems for dealing with big data
- Align the relationships between Hadoop and Spark
- Differentiate between Spark RDDs and Spark Dataframes and when each is appropriate
- locate and explore the Spark.ML documentation
- code along to understand similarities between Pyspark and Python (Pandas/Sklearn)

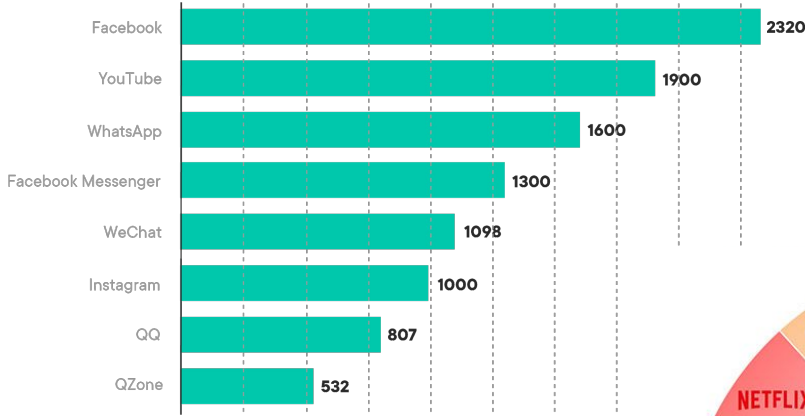
BIG DATA

BIG

makeameme.org

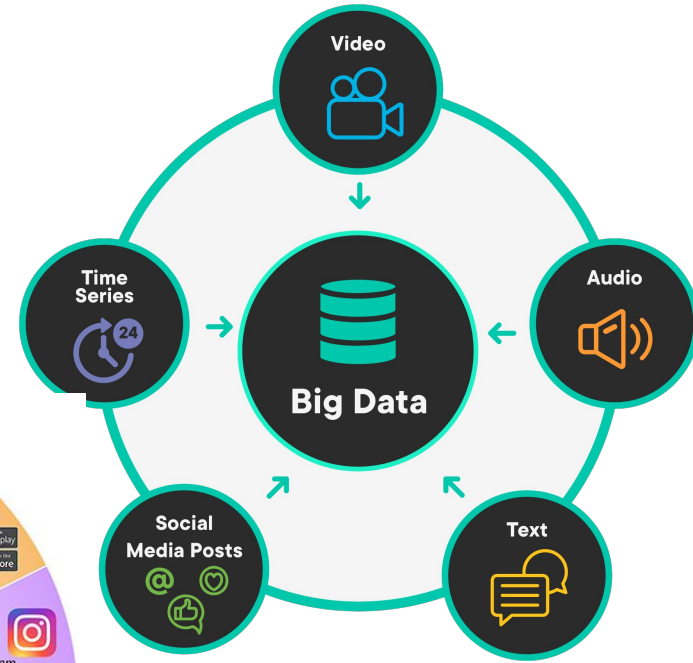
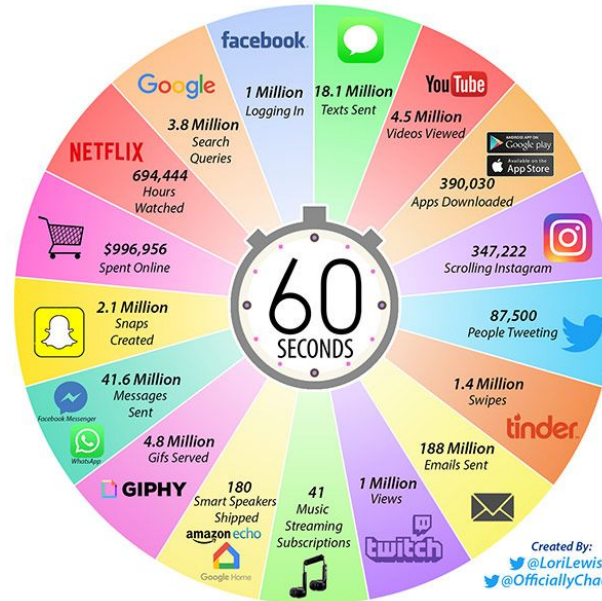


Social Media Networks Ranked by Number of Users (in millions)



BIG VOLUME

HIGH VELOCITY



HIGH VARIETY

Created By:
 @LoriLewis
 @OfficiallyChadd

Where in the Process are we?

1. Business Question
2. Data Science Question(s)
3. Data Acquisition/Cleaning
4. EDA
5. Feature Selection and Engineering
6. Modeling
7. Model Evaluation
8. Addressing Business Questions



Where in the Process are we?

1. Business Question
2. Data Science Question(s)
3. Data Acquisition/Cleaning
4. EDA
5. Feature Selection and Engineering
6. Modeling
7. Model Evaluation
8. Addressing Business Questions



Visualization
& Analytics



Computation



Storage



Distribution &
Data Warehouse



Visualization
& Analytics



Computation



Storage



Distribution &
Data Warehouse



Hadoop 1

- Silos & Largely batch
- Single Processing engine

MapReduce
(Cluster Resource Management
& **Batch** Data Processing)

1 ° ° ° °
HDFS
(Hadoop Distributed File System)

Hadoop 2 w/YARN

- Multiple Engines, Single Data Set
- Batch, Interactive & Real-Time

Batch
MapReduce

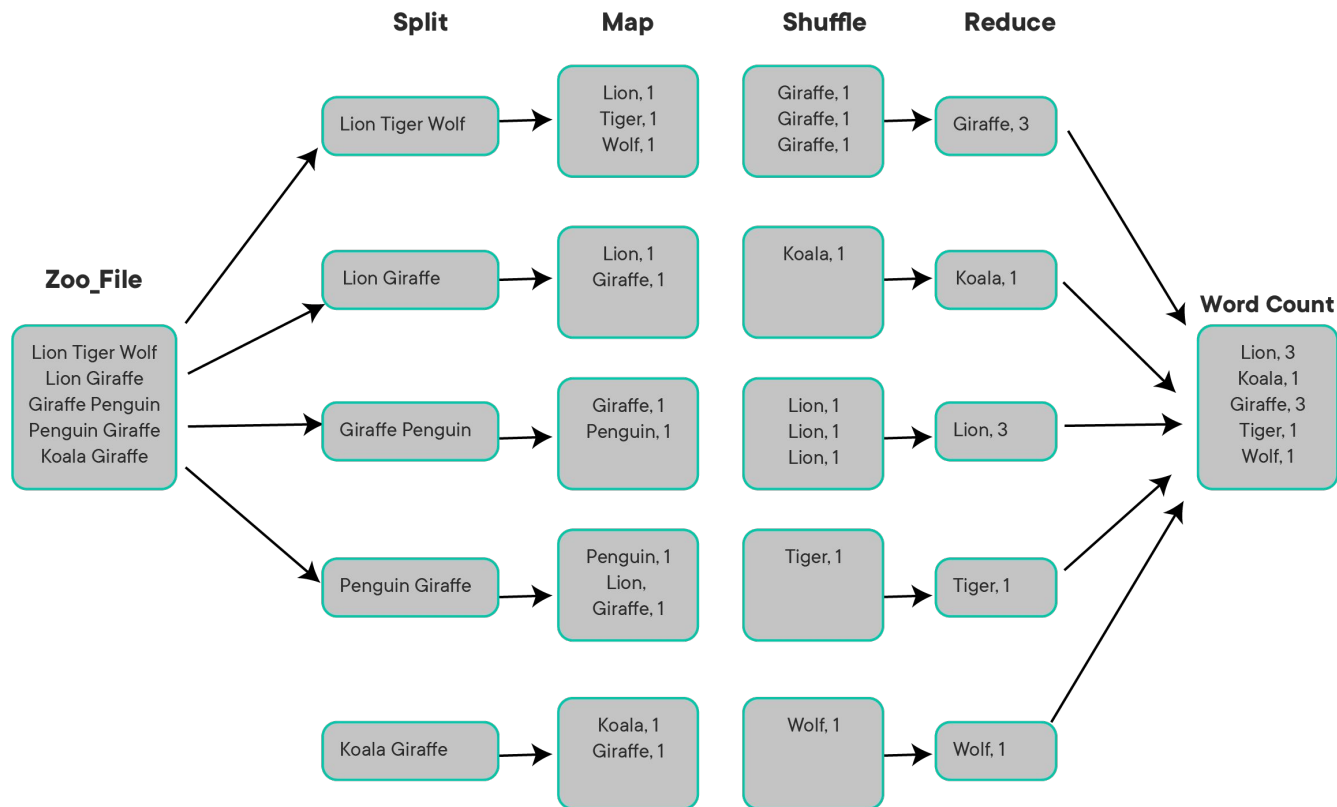
Interactive
Others

Real-Time
Others

YARN: Data Operating System
(Cluster Resource Management)

1 ° ° ° ° ° °
HDFS
(Hadoop Distributed File System) N

MapReduce



Applications Run Natively IN Hadoop

Pig

Script

Hive

SQL

HBase

NoSQL

Accumulo

NoSQL

Storm

Stream

Solr

Search

Spark

In-Memory

Cascading

Java

Others

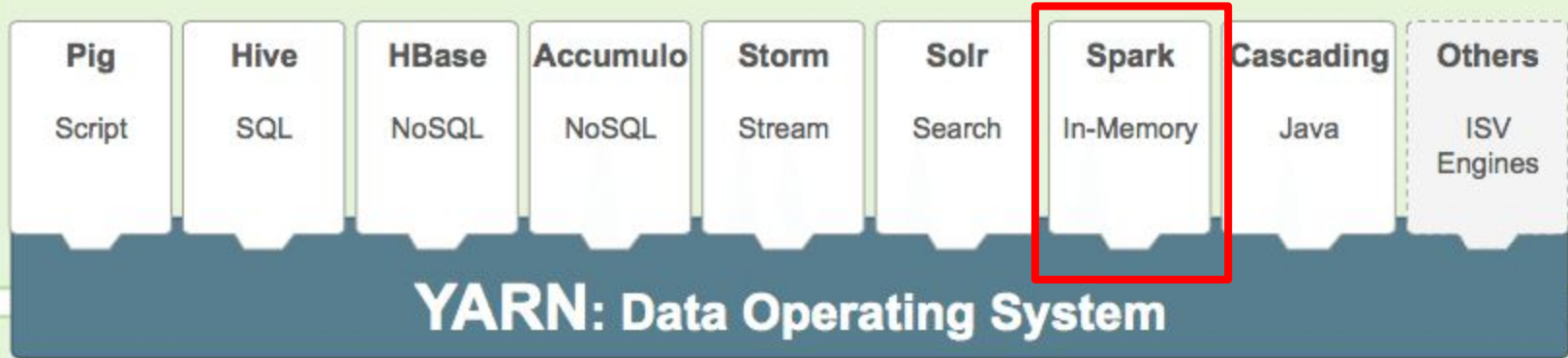
ISV
Engines

YARN: Data Operating System

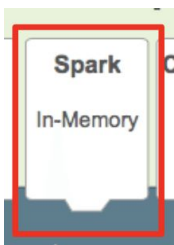
HDFS

(Hadoop Distributed File System)

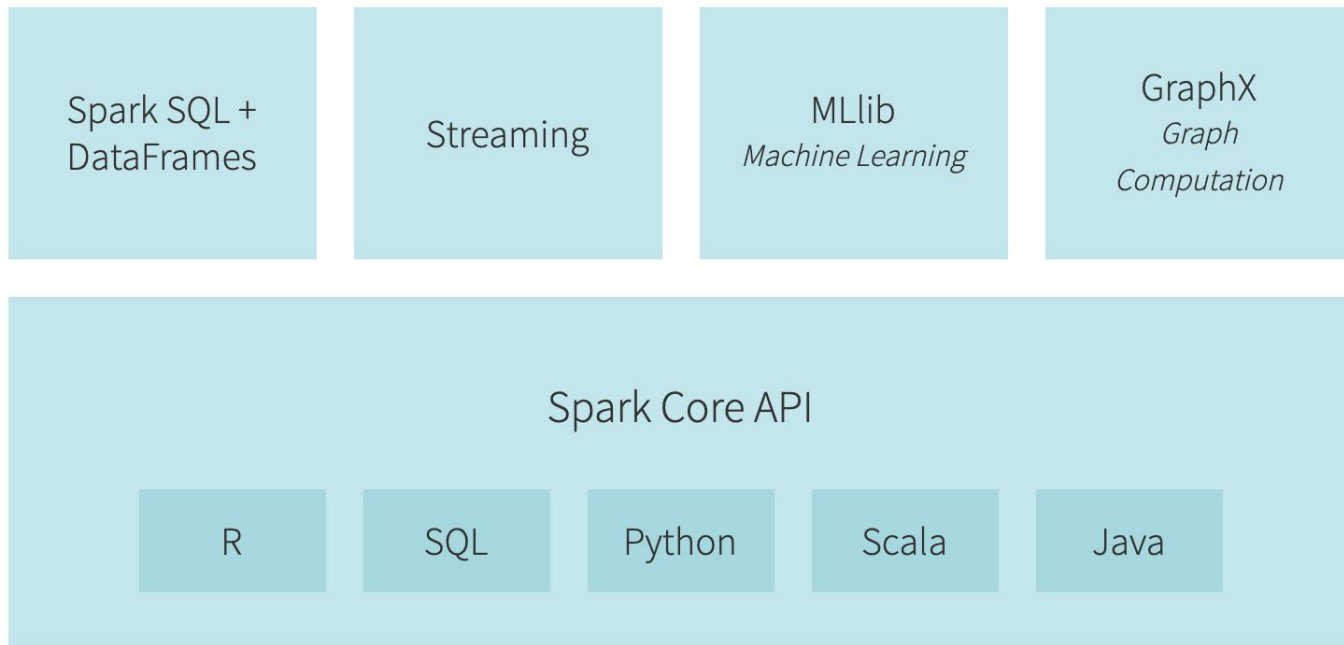
Applications Run Natively IN Hadoop



HDFS
(Hadoop Distributed File System)



Apache Spark Ecosystem



Reminder:

Dictionary



API

noun

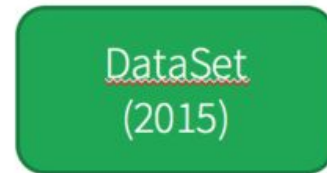
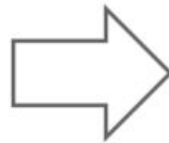
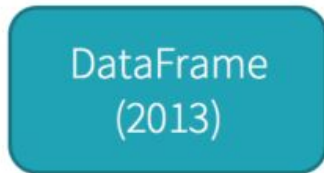
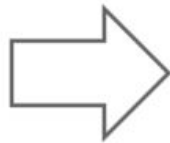
COMPUTING

a set of functions and procedures allowing the creation of applications that access the features or data of an operating system, application, or other service.



Translations, word origin, and more definitions

Spark Data Objects



Distribute collection
of JVM objects

Functional Operators (map,
filter, etc.)

Distribute collection
of Row objects

Expression-based operations
and UDFs

Logical plans and optimizer

Fast/efficient internal
representations

Internally rows, externally
JVM objects

Almost the “Best of both
worlds”: **type safe + fast**

But slower than DF
Not as good for interactive
analysis, especially Python

Space Efficiency

Memory Usage when Caching



Space Efficiency

Memory Usage when Caching

Datasets

- high-level expressions, filters, maps, aggregation, averages, sum
- SQL queries, columnar access and use of lambda functions

RDDs

- Unstructured data
- you don't care about imposing a schema, such as columnar format, while processing or accessing data attributes by name or column

0

15

30

45

60

Data Size (GB)

Space Efficiency

Memory Usage when Caching



Datasets

- high-level expressions, filters, maps, aggregation, averages, sum
- SQL queries, columnar access and use of lambda functions

RDDs

- Unstructured data
- you don't care about imposing a schema, such as columnar format, while processing or accessing data attributes by name or column

0

15

30

45

60

Data Size (GB)

DataFrame or RDD?

1. You are grabbing live tweets about the the 2020 Political Election.
2. You have an RDD of data that you wish to use to build a predictive model.
Should you leave it as an RDD or transform it to a DataFrame?
3. You want to access audio and video stored on your HDFS

Lets Code