# Introduction to Data Engineering and Spark Architecture

# Agenda

- **Peek into Data Engineering & Big Data**

    - **Data Engineering**

    - **Big Data Frameworks**

- **Introduction to Spark**

    - **Spark vs. Hadoop**

    - **Spark Architecture**

- **Setting up a Spark Cluster**

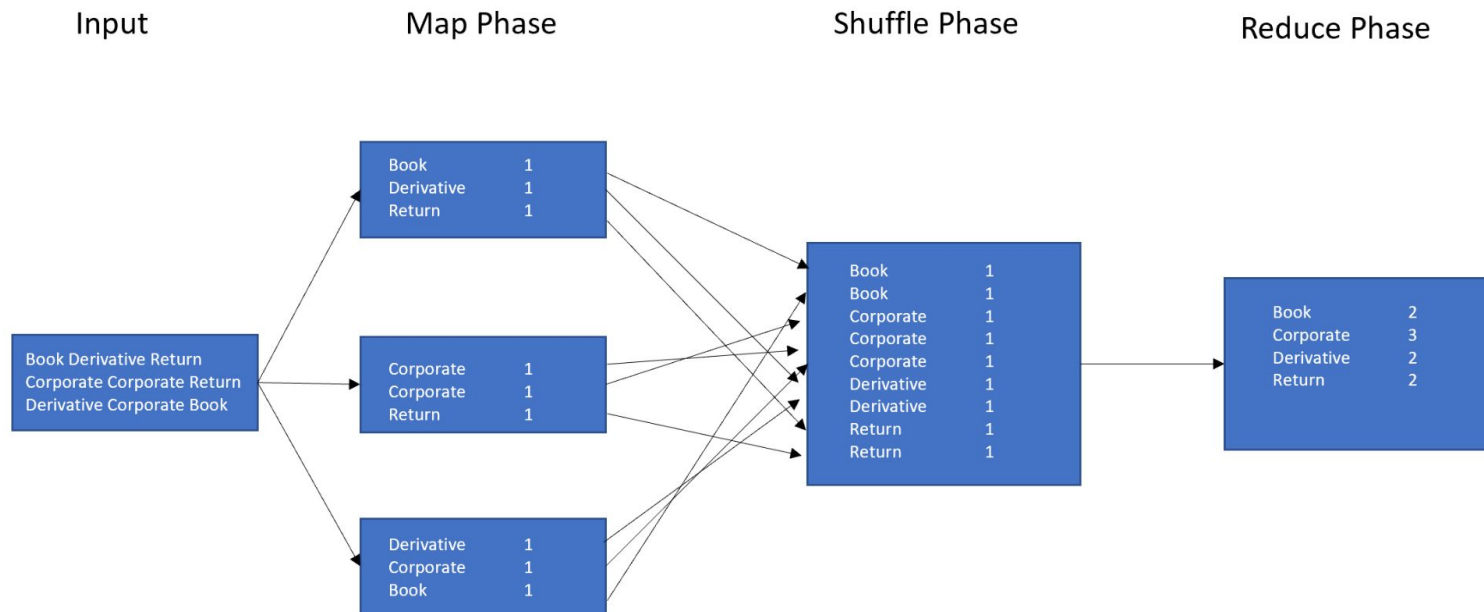# Peek into Data Engineering & Big Data

# Data Engineering

- Data Engineers support Data Scientists.
  - In charge of designing, creating, deploying, and supporting data pipelines.
  - Depending of the side of the company, the Data Scientist can be a "Full Stack Data Scientist": in charge of its own Data Engineering.
- Skills needed
  - Computer Science
  - Business knowledge
  - Database Expertise (SQL and NoSQL)
  - Big Data Architectures
- Why Become a Data Engineer?
  - Explosion on roles over the last 5 years.
  - Unlike Data Scientist, there is no clear path.
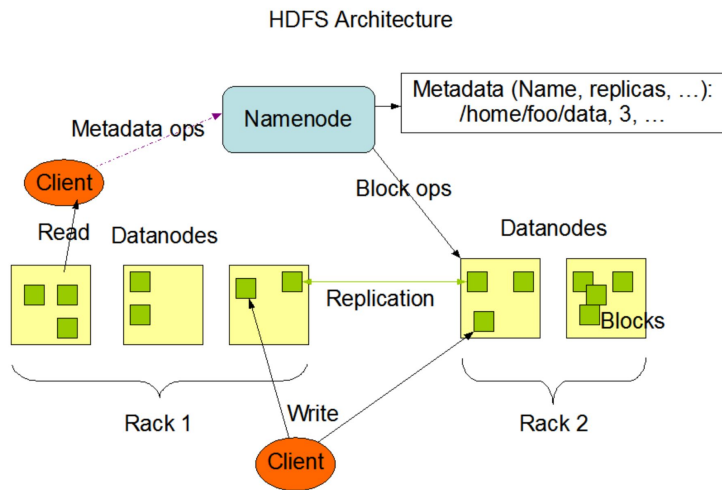
# Big Data Frameworks

- Big Data
  - When is considered big? Loosely defined.
  - Depends on both external and internal factors.
  - Rule of Thumb: "*Too Big to fit in a Pandas Dataframe"*.
- Parallel computation
  - Divide and conquer - Map Reduce!
  - Map phase - Activities that can be done independently.
  - Reduce phase - Aggregation done at the end.
  - Data structure: Key-Value pairs

# Big Data Frameworks

● MapReduce example

| Input | Map Phase | Shuffle Phase | Reduce Phase |
|---|---|---|---|

**Input:**
Book Derivative Return
Corporate Corporate Return
Derivative Corporate Book

**Map Phase:**

Book 1
Derivative 1
Return 1

Corporate 1
Corporate 1
Return 1

Derivative 1
Corporate 1
Book 1

**Shuffle Phase:**

Book 1
Book 1
Corporate 1
Corporate 1
Corporate 1
Derivative 1
Derivative 1
Return 1
Return 1

**Reduce Phase:**

Book 2
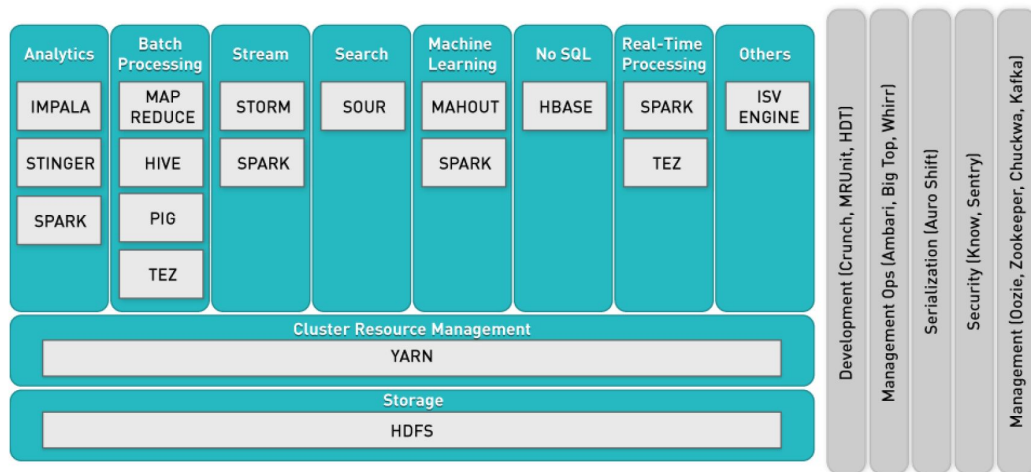Corporate 3
Derivative 2
Return 2

# Big Data Frameworks



HDFS Architecture

- Distributed computing
  - Instead of a big powerful machine - Several simpler ones
  - Code lives in Master - Work is done in Workers.

# Big Data Frameworks

**Apache Hadoop 2.0 Ecosystem**

| Analytics | Batch Processing | Stream | Search | Machine Learning | No SQL | Real-Time Processing | Others | Development (Crunch, MRUnit, HDT) | Management Ops (Ambari, Big Top, Whirr) | Serialization (Auro Shift) | Security (Know, Sentry) | Management (Oozie, Zookeeper, Chuckwa, Kafka) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| IMPALA | MAP REDUCE | STORM | SOUR | MAHOUT | HBASE | SPARK | ISV ENGINE | | | | | |
| STINGER | HIVE | SPARK | | SPARK | | TEZ | | | | | | |
| SPARK | PIG | | | | | | | | | | | |
| | TEZ | | | | | | | | | | | |

**Cluster Resource Management**

YARN

**Storage**

HDFS

http://incubator.apache.org/projects/

- Hadoop
  - Distributed File System (HDFS)
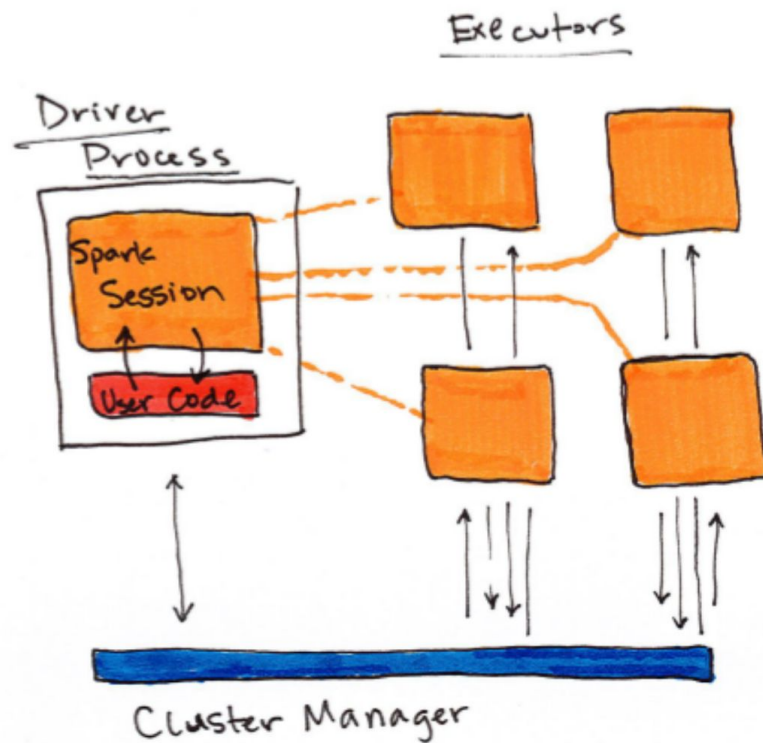  - Manages both work distribution and fault tolerance.

# Introduction to Spark

# Spark vs. Hadoop

- MapReduce has been the major framework for distributed computing
  - Hadoop's limitations include programmability and performance.
  - Computational frameworks are becoming specialized.
- Spark: 100x faster than Hadoop
  - Spark is the Compute Engine - Hadoop still provides the environment.
- Specialized libraries for machine learning, graph processing, and database management.
- APIs include:
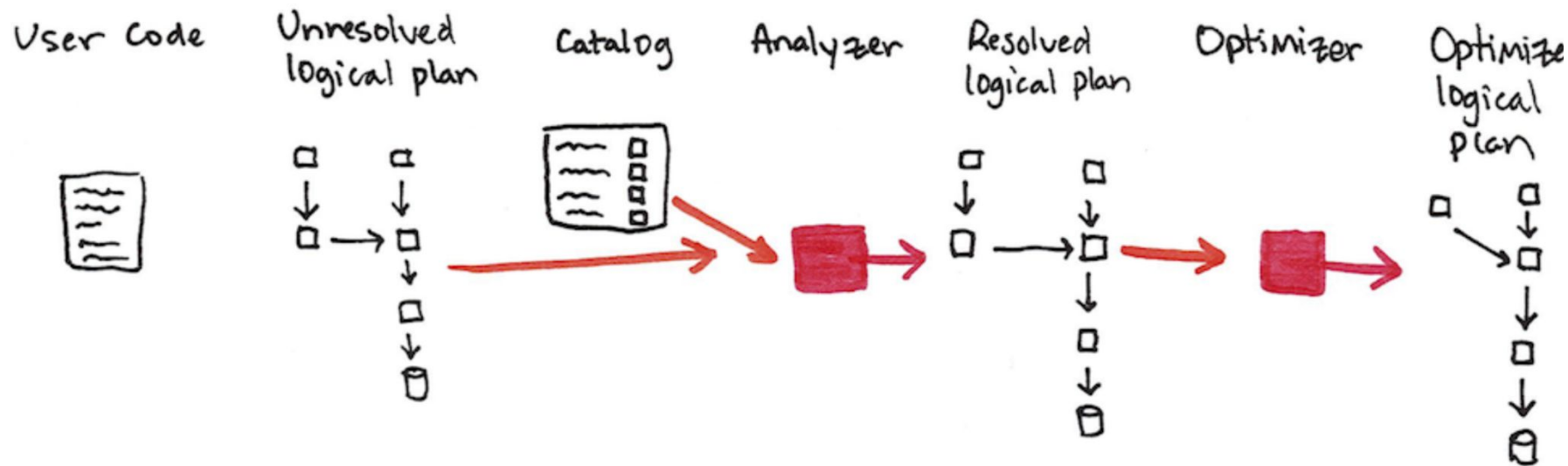  - Java
  - Scala
  - Python
  - R

# Spark Architecture (I)

- Driver and Executors

# Spark Architecture (II)

- SparkContext and Lazy Evaluation

**Happy Learning !**