

CSCE 5300: Introduction to Big Data and Data Science

Lesson 1

Overview

Overview

- Evaluation Criteria
- Topics to be covered
- Installations
- In class Exercise

Grading Criteria

- . 20% Quizzes (individual)
- . 25% In-class Tasks
- . 30% Project
- . 25% Exam

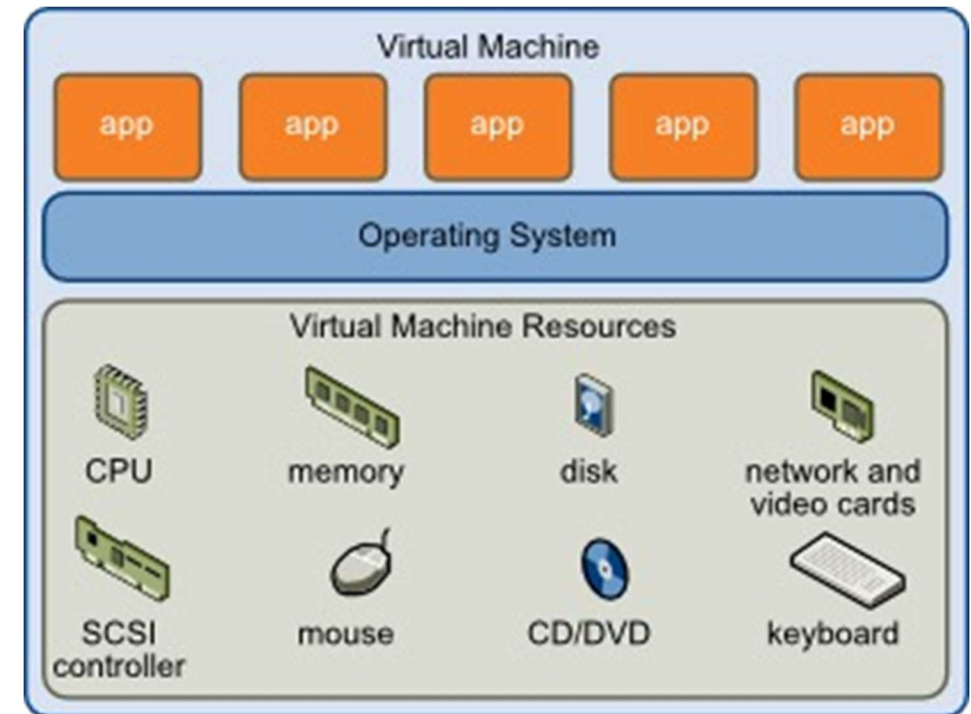
Topics to be Covered

- Big Data Overview, Installations
- HDFS / Map Reduce / Big Data Applications
- Hadoop Dependent Query Based No SQL Database Hive
- Hadoop to SQL Parallel Transfer Engine: Sqoop
- Parallel Indexing: Solr & Lucene
- Independent Column Based No SQL Database: Cassandra
- Spark Programming with RDDs and applications
- Spark: Data Frames and SQL
- Machine Learning and Big Data Analytics Applications
- Data Visualization, Deep Learning Concepts
- Spark with RDD and streaming
- GraphX, GraphFrames, Graph Analytics Applications
- Parallel Computing

Cloudera

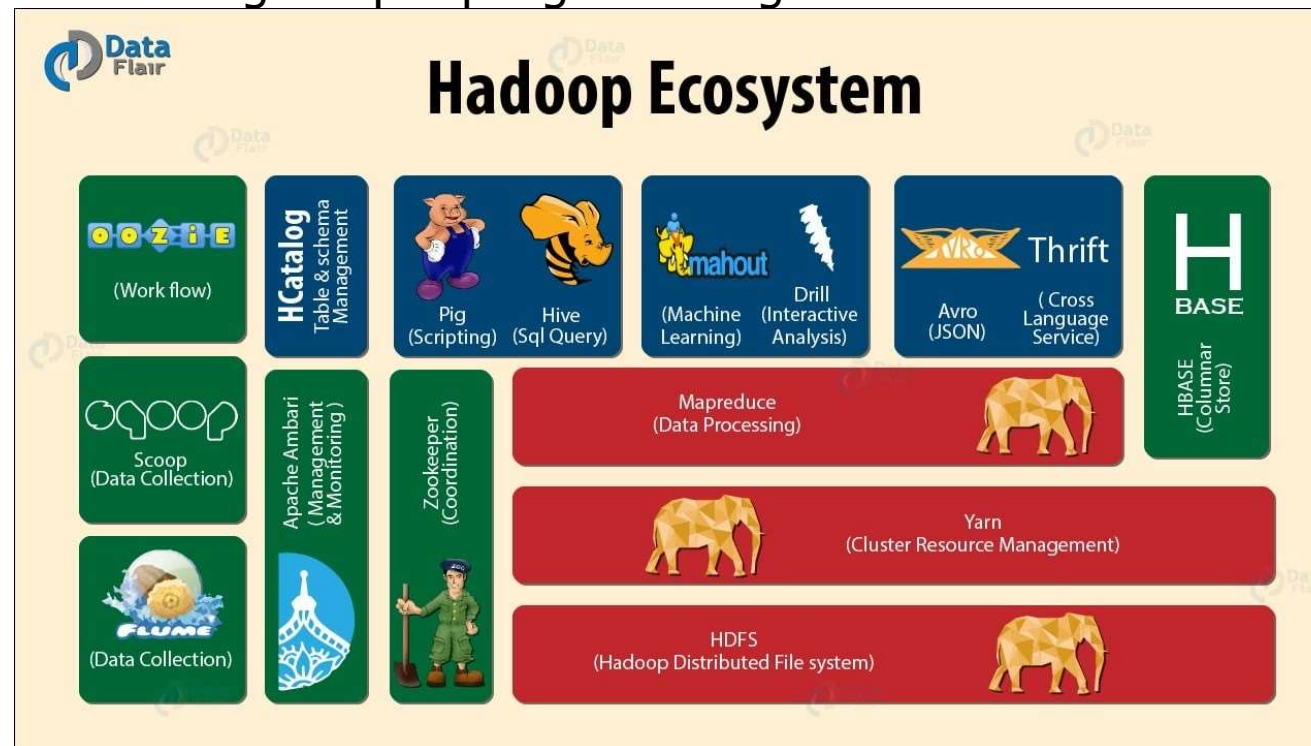
Virtual Machine

- In computing, a **virtual machine (VM)** is an emulation of a computer system
- Virtual machines are based on computer architectures and provide functionality of a physical computer.
- Their implementations may involve specialized hardware, software, or a combination



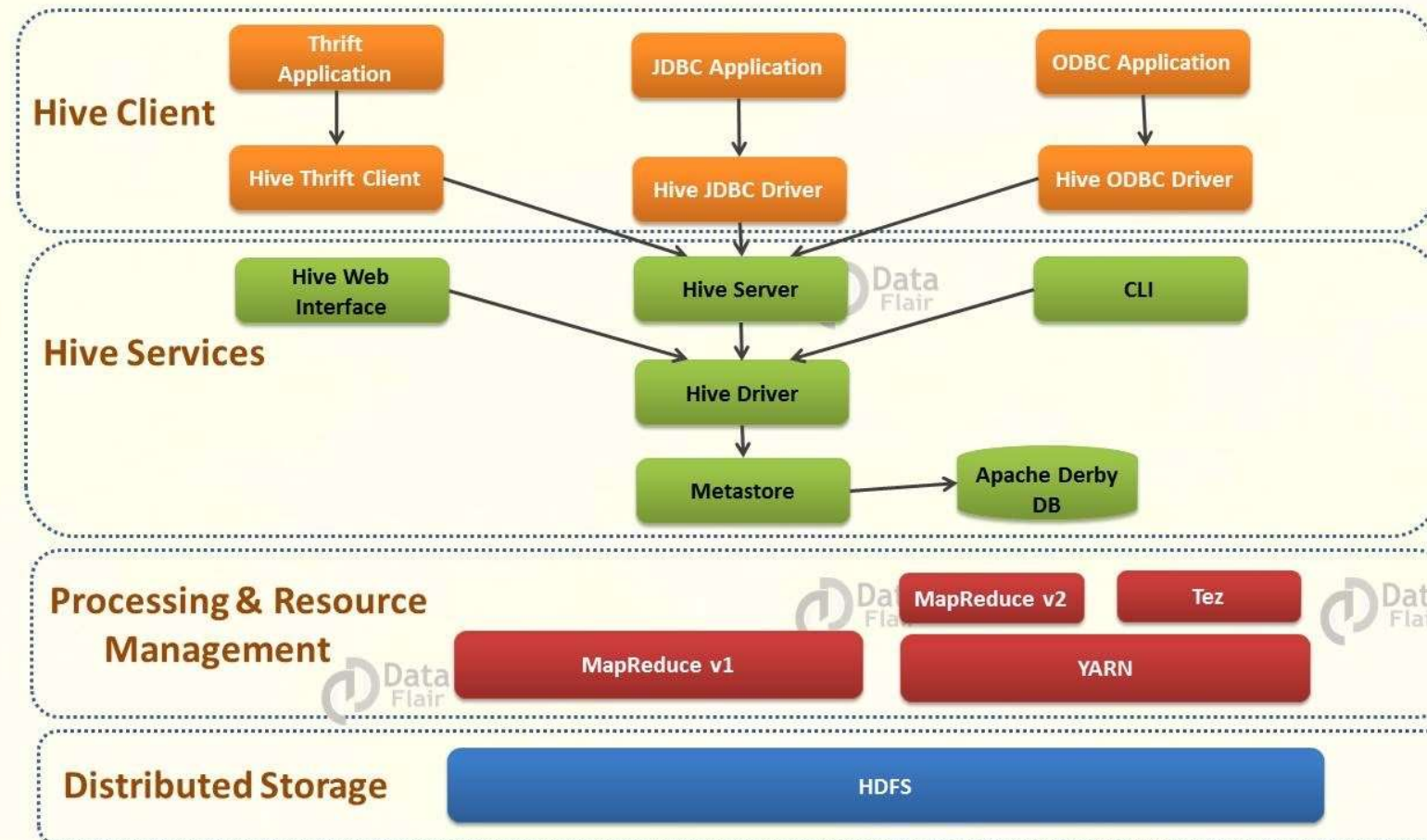
Hadoop Eco-system

A framework that allows for the distributed processing of large data sets across clusters of computers using simple programming models



Source: <https://data-flair.training/blogs/hadoop-ecosystem-components/>

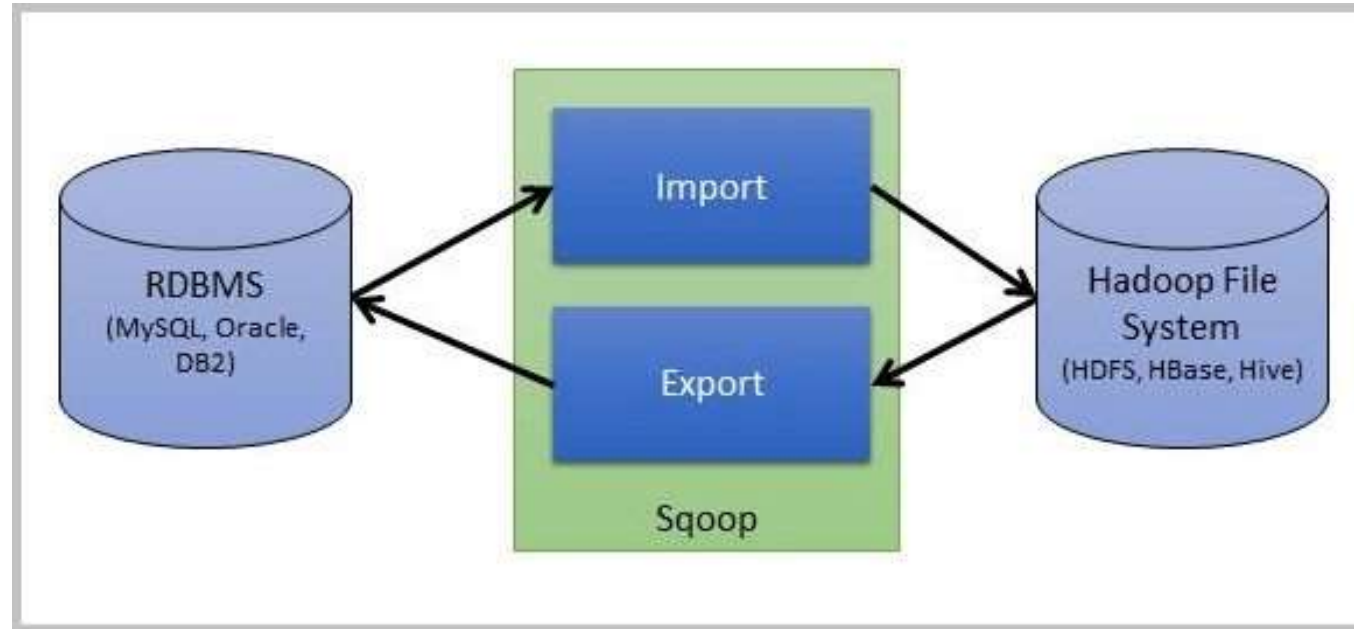
Hive Architecture & its Components



Source: <https://data-flair.training/blogs/apache-hive-architecture/>

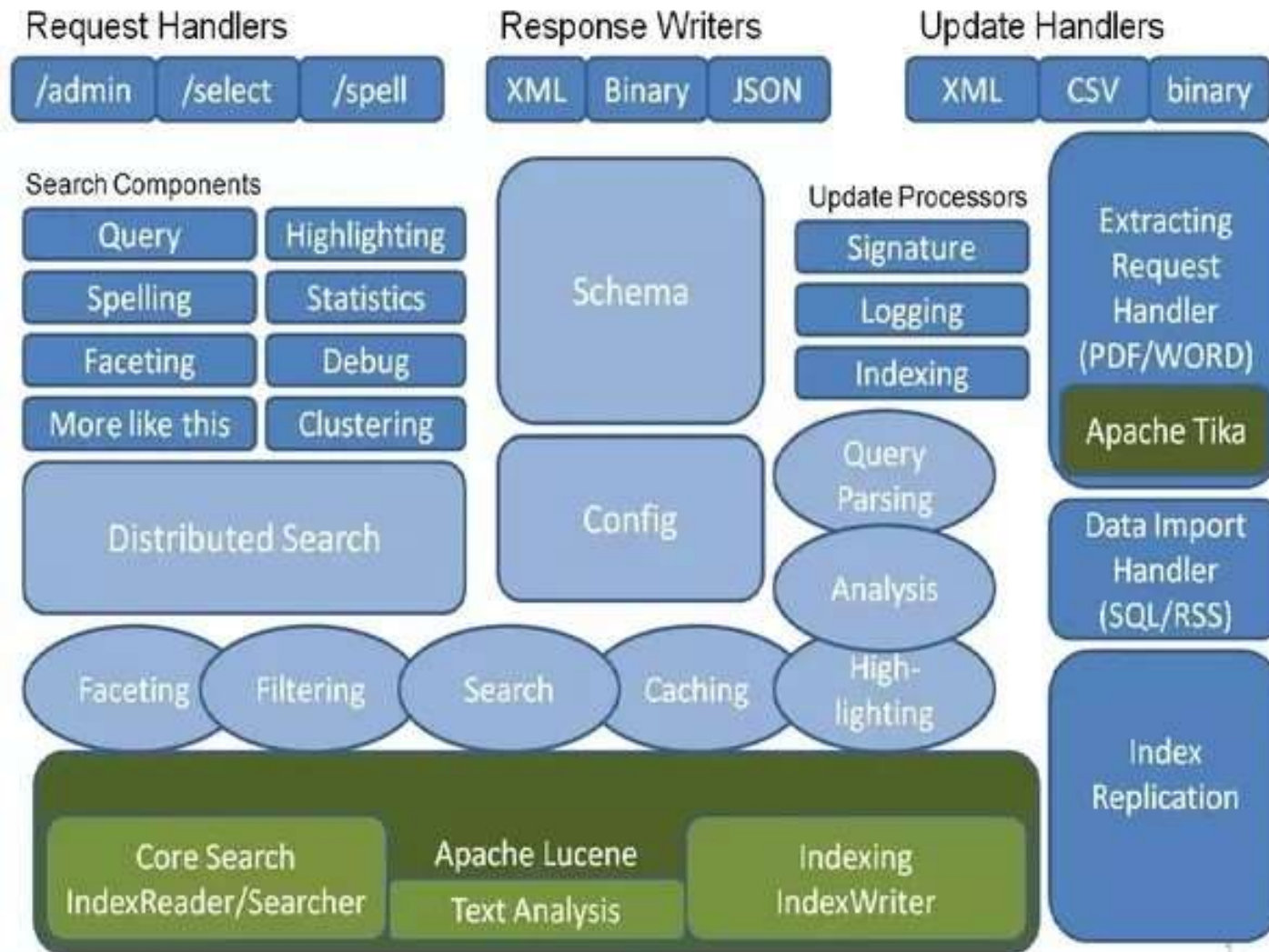
Sqoop

Application for transferring data between relational databases and Hadoop



Source: <https://www.hdfstutorial.com/sqoop-architecture/>

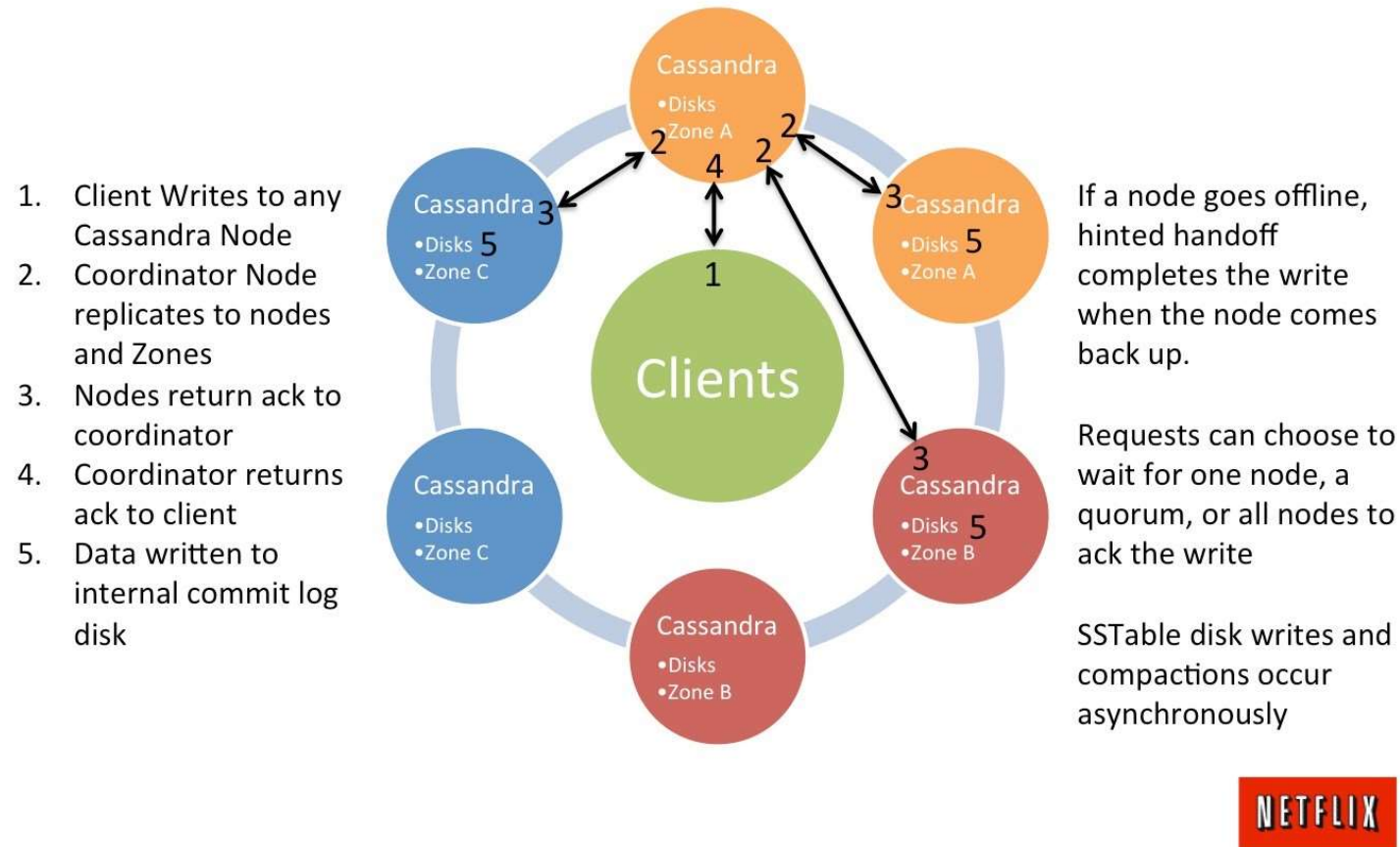
Lucene/Solr Architecture



Source: <https://www.quora.com/What-is-the-internal-architecture-of-Apache-solr>

Cassandra Write Data Flows

Single Region, Multiple Availability Zone



Source: <https://intellipaat.com/tutorial/cassandra-tutorial/brief-architecture-of-cassandra/>