# Popular Open Source Data Processing Frameworks

Giri Iyengar

Cornell University

*gi43@cornell.edu*

April 23, 2018

# Agenda for the week

- Pig
- Spark
- Storm
- BlinkDB
- Druid
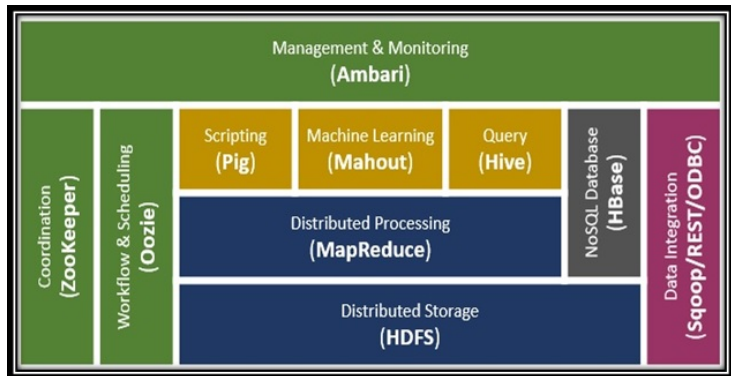
# Overview

1. Pig

2. Spark
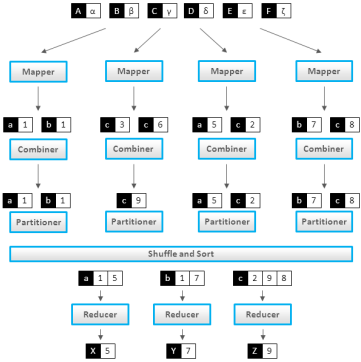
# Hadoop



Figure: Hadoop Ecosystem
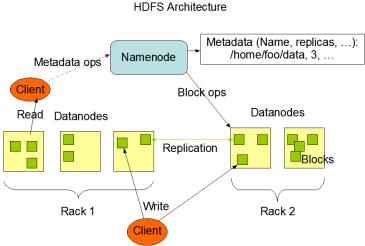
# Hadoop



Figure: MR Framework



Figure: HDFS

# Pig Demo



Figure: Pig Reference Manual

# Hadoop

- Use commodity hardware to achieve super stable, reliable, data processing
- Reliable data storage via replication - HDFS
- Generic Computation approach - Map Reduce
- Extremely well-suited for batch processing on 1000s of nodes handling Petabytes of data

# Since the advent of Hadoop

- Speed and sophistication required for data processing has grown tremendously
- Complex algorithms like Machine Learning and Graph Analysis are much more common
- E.g. ML requires multiple passes over the data – not suited for Map Reduce style computing
- Streaming analysis of real-time data is increasingly important
- Both one-pass aggregations and multi-pass analysis applications need to be supported

# Overview

# Spark

## Apache Spark

Created by Matei Zaharia as part of his Doctoral work at UC Berkeley.
Designed to address some of the limitations observed with Hadoop.

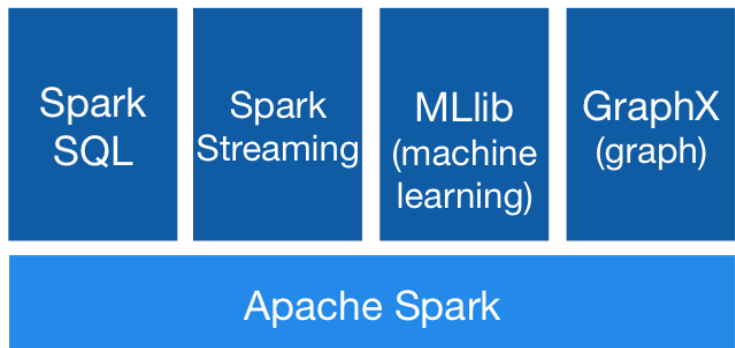Stated goal is to scale to 10s of thousands of compute nodes

# Spark Stack



Figure: Spark Stack

# Spark Cluster Overview



Figure: Spark Stack

# Spark RDD

## Resilient Distributed Dataset

Spark revolves around the concept of a resilient distributed dataset (RDD), which is a fault-tolerant collection of elements that can be operated on in parallel. There are two ways to create RDDs: parallelizing an existing collection in your driver program, or referencing a dataset in an external storage system, such as a shared filesystem, HDFS, HBase, or any data source offering a Hadoop InputFormat.

# RDD

- Basic Abstraction in Spark
- Immutable, Partitioned collection of elements that can be operated in parallel
- Supports Lazy evaluations

# Map Reduce Intermediate Data

# Spark Intermediate Data

# Properties of RDD

- In-Memory
- Lazy
- Fault-Tolerant
- Immutability
- Partitioned
- Persistent
- Parallel

- Location-Stickiness
- Typed
- Coarse-Grained Operations (whole RDD and not individual elements)
- No limitations (bound by available system memory)

# Spark DataFrames

## Spark DataFrames

This API is inspired by data frames in R and Python (Pandas), but designed from the ground-up to support modern big data and data science applications. It is an extension to the existing RDD API.

# Spark DataFrames

- Ability to scale from kilobytes of data on a single laptop to petabytes on a large cluster Support for a wide array of data formats and storage systems
- State-of-the-art optimization and code generation through the Spark SQL Catalyst optimizer
- Seamless integration with all big data tooling and infrastructure via Spark
- APIs for Python, Java, Scala, and R