# Hadoop Components Overview

# Contents

- ➢ Hadoop Components

- ➢ Hadoop Common

- ➢ HDFS

- ➢ MapReduce

- ➢ Hadoop Tools

**Components of Hadoop:**

- ➢ Hadoop Common

- ➢ HDFS

- ➢ Yarn

- ➢ Map Reduce

➢ A set of utilities that supports the Hadoop subprojects

➢ Provides access to the file systems supported by Hadoop

➢ Hadoop Common package contains:

  ✓ The JAR files and scripts necessary to start Hadoop

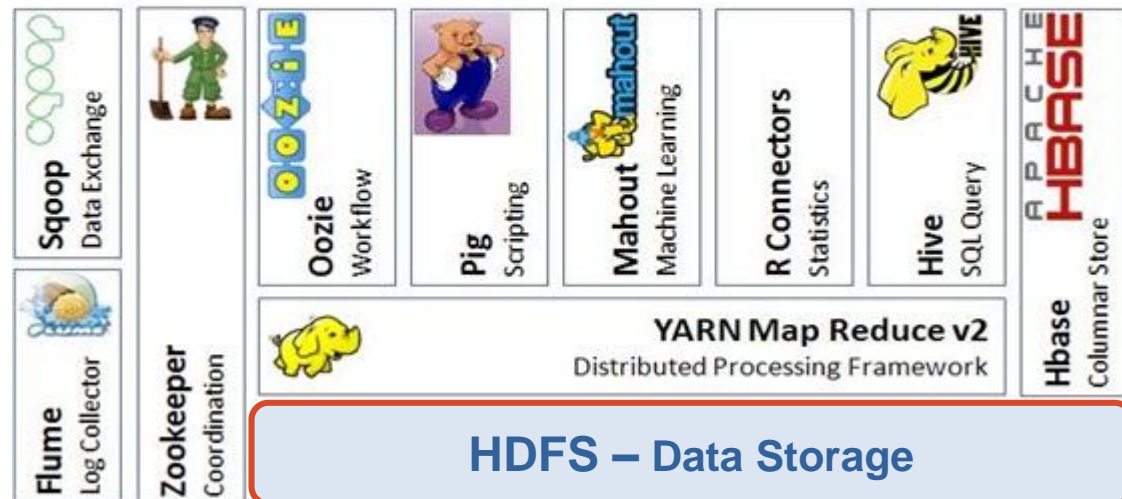  ✓ Source code, documentation and a contribution section with projects from the Hadoop Community

HDFS - **H**adoop **D**istributed **F**ile **S**ystem

- ➢ Primary storage system for Hadoop

- ➢ Distributed, portable, scalable file system written in Java

- ➢ Files are broken down and stored in multiple machines

- ➢ Designed for large scale distributed data processing

- ➢ Follows master/slave architecture

Apache Hadoop Ecosystem

**HDFS is best suited for:**

- ➢ Highly fault-tolerant

- ➢ Designed to deploy on low-cost hardware

- ➢ Suitable for applications that have large datasets(Giga to Terabytes)

- ➢ Enables streaming access to file system data

**HDFS is not good for:**

- ➢ Low Latency Data Access

- ➢ Lot of Small Files

- ➢ Multiple Writers, Arbitrary File Modifications

- ➢ Programming framework (library and runtime) for analyzing data sets stored in HDFS
- ➢ MapReduce jobs are composed of two functions:

$$\texttt{map()} \rightarrow \texttt{reduce()}$$

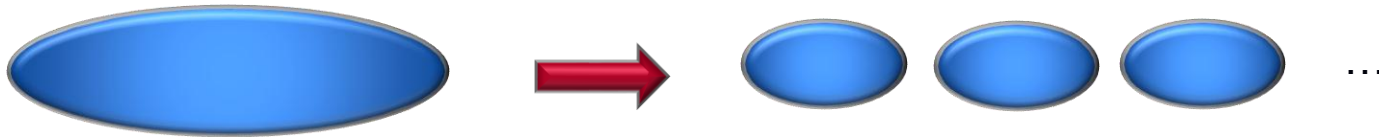*sub-divide & conquer*     *combine & reduce cardinality*

- ➢ User only writes the Map and Reduce functions

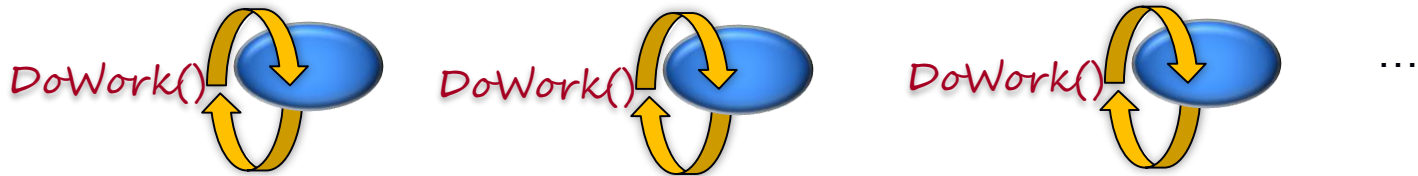- ➢ MR framework provides all the "glue" and coordinates the execution of the Map and Reduce jobs on the cluster
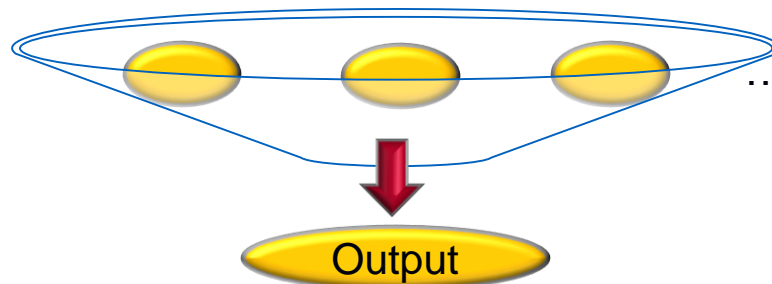
**MAP**

**REDUCE**

Essentially, it's…

1. Take a large problem and divide it into sub-problems
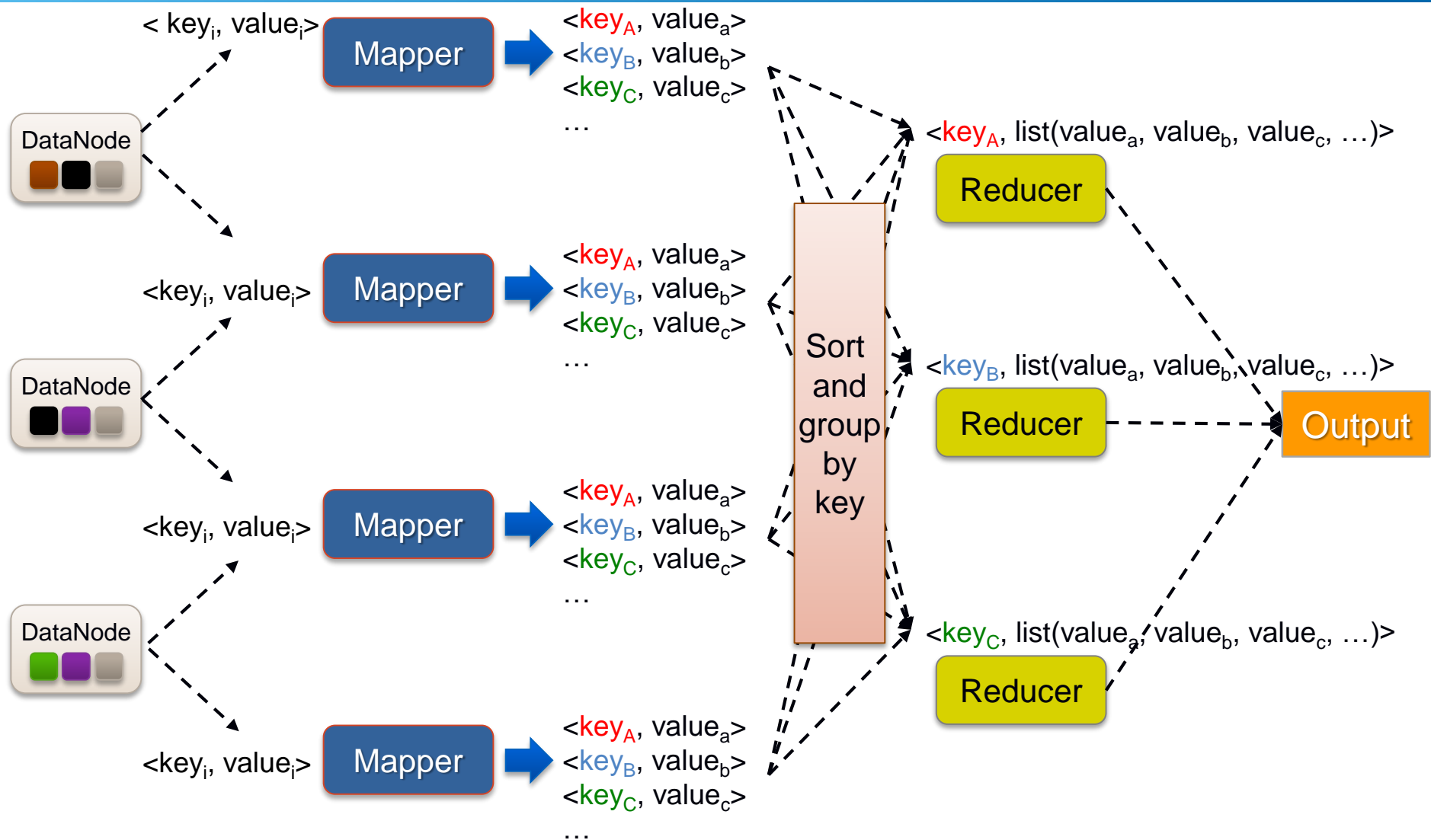
2. Perform the same function on all sub-problems

   DoWork()    DoWork()    DoWork()    …

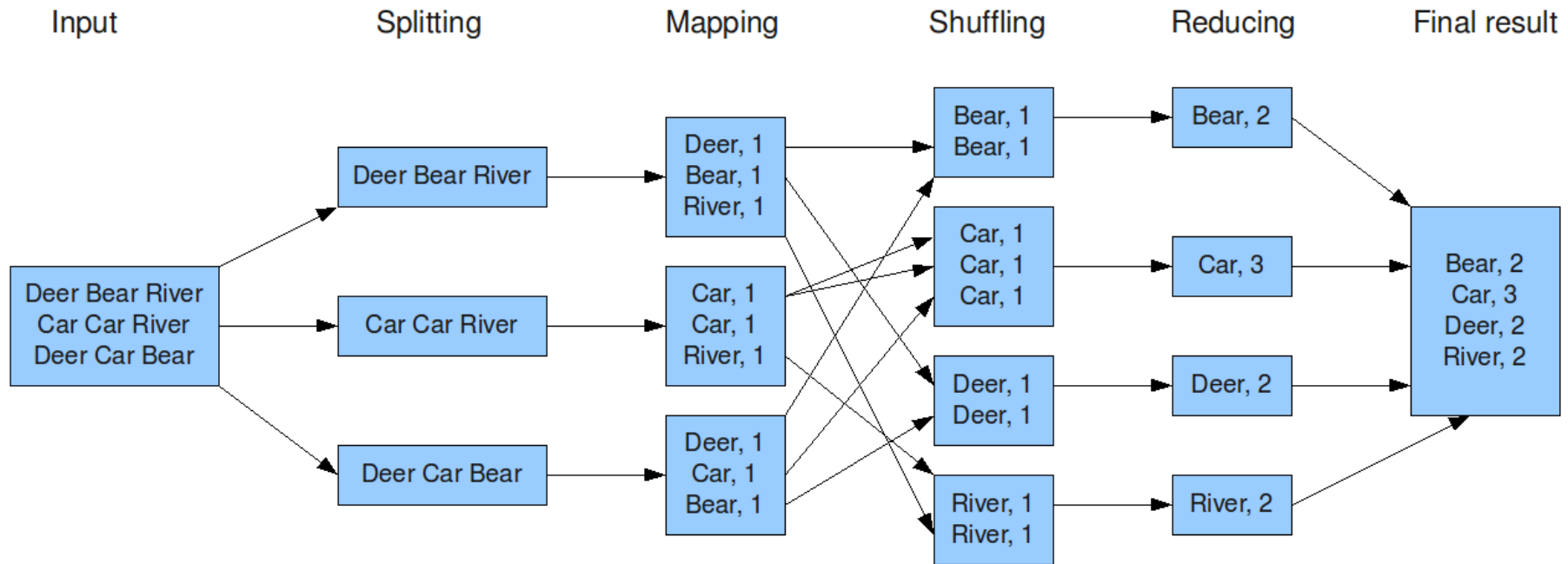3. Combine the output from all sub-problems

   Output

A key-value pair (KVP) is a set of two linked data items:

Key - A unique identifier for data item

Value – Either the data that is identified or a pointer to the location of that data
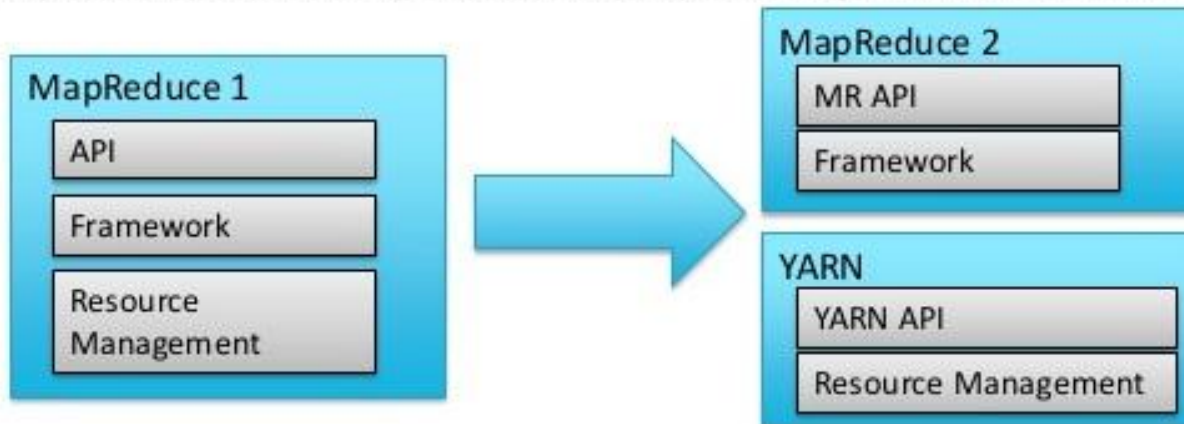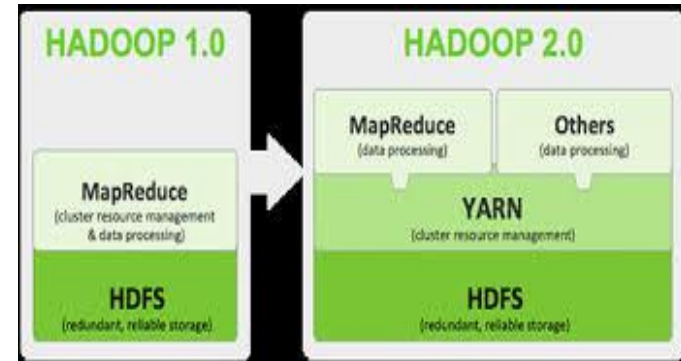
The overall MapReduce word count process

## Next generation MapReduce 2



- **MapReduce 1 ("Classic") has three main components**
  - API – for user-level programming of MR applications
  - Framework – runtime services for running Map and Reduce processes, shuffling and sorting, etc.
  - Resource management – infrastructure to monitor nodes, allocate resources, and schedule jobs

- **MapReduce 2 ("NextGen") moves Resource Management into YARN**

# Current MapReduce vs YARN

## Hadoop MapReduce

**Job Tracker (Master)**
- Resource management
- Job lifecycle management
  - Scheduling, progress monitoring, fault tolerance

**Task Tracker (per node)**
- Launch tasks
- Report status to Job tracker

## Hadoop YARN

**Resource Manager (Master)**
- Resource management
- Scheduling

**Application Master (per app)**
- Job lifecycle management

**Node Manager (per node)**
- Launch Containers
- Monitor resource usage
- Report to RM

MapReduce itself is an Application on YARN

## Pig

➢ Pig is a platform for analysing large data sets that consists of a high-level language for expressing data analysis programs

➢ Consists of :

✓ PigLatin : The high-level language

✓ A Run-time environment where PigLatin programs are executed

## HBase

➢ Column-oriented database management system on top of HDFS

➢ HBase system comprises a set of tables

➢ Allows Attributes (Columns) to be grouped together into "column families"

➢ All the elements of a column family are all stored together

TATA CONSULTANCY SERVICES

## Hive

➢ A data warehouse infrastructure built on top of Hadoop

➢ Provides tools to enable :

  ✓ Easy data summarization

  ✓ Ad hoc querying

  ✓ Analysis of large datasets stored in Hadoop files

## Sqoop - SQL To Hadoop

➢ Sqoop is a command-line tool with the following capabilities:

  ✓ Imports tables or entire databases to files in HDFS

  ✓ Generates Java classes to allow you to interact with imported data

  ✓ Imports from SQL databases straight to Hive data warehouse

## ZooKeeper



➢ A high-performance coordination service for distributed applications

➢ An open source Apache project

➢ Provides a centralized infrastructure and services that enable synchronization across a cluster

➢ Maintains common objects needed in large cluster environments

➢ Each client machine communicates with one of the many ZooKeeper servers in a cluster to retrieve, update its synchronization information

TATA
CONSULTANCY
SERVICES

## Oozie

➢ Is a workflow/coordination service to manage data processing jobs for Apache Hadoop

➢ Supports all types of Hadoop jobs and is integrated with the Hadoop stack

➢ Users can specify execution frequency and can wait for data arrival to trigger an action in the workflow

## Flume

➢ Service for efficiently moving large amounts of data soon after the data is produced

➢ It is centrally managed & allows for intelligent dynamic management

➢ Its main goal is to deliver data from applications to Hadoop's HDFS

## Mahout

➢ A scalable machine learning and data mining library

➢ An Apache project to produce free implementations of machine learning algorithms on the Hadoop platform

- *Understanding Big Data*- Chris Eaton, Dirk Deroos, Tom Deutsch, George Lapis, Paul Zikopaulos

- *Hadoop: The Definitive Guide*-O'Reilly

- *Hadoop In Action*-by Chuck Lam

- archive.cloudera.com/cdh3/hue

- github.com/nathanmarz/cascalog

- archive.cloudera.com/cdh3/hue

- http://en.wikipedia.org/wiki/Pentaho

- http://www.pentaho.com/hadoop/

- http://katta.sourceforge.net/

- http://developer.yahoo.com/hadoop/tutorial/module4.html

# Thank You