# The Hadoop Ecosystem:

# So much free stuff!

# After this video you will be able to..

- Differentiate the major layers in the Hadoop ecosystem

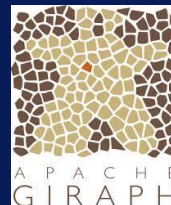- Recognize key tools of the Hadoop ecosystem including HDFS, YARN, and MapReduce

Now there's over a 100!

# One possible layer diagram for Hadoop

# One possible layer diagram for Hadoop



Higher levels:
Interactivity

Zookeeper

Hive

Pig

MapReduce

Giraph

Storm

Spark

Flink

YARN

HDFS

HBase

Cassandra

MongoDB

Lower levels:
Storage and scheduling

Distributed file system as foundation

Scalable storage

Fault tolerance

Zookeeper | Hive | Pig | Giraph | Storm | Spark | Flink | MapReduce | YARN | HDFS | HBase | Cassandra | MongoDB

# Flexible scheduling and resource management

YARN schedules jobs on > 40,000 servers at Yahoo

Zookeeper

Hive

MapReduce

Giraph

Storm

Spark

Flink

HBase

Cassandra

MongoDB

YARN

HDFS

# Simplified programming model

Map → apply()

Reduce → summarize()



Zookeeper

Hive

Pig

MapReduce

Giraph

Storm

Spark

Flink

Base

dra

OB

Google used MapReduce
for indexing web sites

# Higher-level programming models

## Pig = dataflow scripting

## Hive = SQL-like queries

Zookeeper

Hive

Pig

Giraph

Storm

Spark

Flink

MapReduce

HDFS

Pig created at Yahoo,
Hive created at Facebook

Real-time and in-memory processing

In-memory → 100x faster for some tasks

Zookeeper | Hive | Pig | Giraph | Storm | Spark | Flink | HBase | Cassandra | MongoDB

MapReduce

YARN — hadoop

HDFS — hadoop HDFS

# NoSQL for non-files

## Key-values

## Sparse tables



Hive

Pig

Giraph

Storm

Spark

Flink

MapReduce

HBase

Cassandra

MongoDB

HBase used for Facebook's Messaging Platform

HDFS

# Zookeeper for management

## Synchronization

## Configuration

## High-availability

Zookeeper

Hive

Pig

Spark

cassandra

Ma

Cas

Mongo

Created by Yahoo to wrangle services named after animals

HDFS    hadoop HDFS

All these tools are open-source

# All these tools are open-source



## Large community for support

All these tools are open-source

Large community
for support

Download separately
or part of pre-built image

Growing number of open-source tools