# The Hadoop Ecosystem



## Hadoop Overview

### What is Hadoop?

### Why Hadoop?

### Hadoop History

## World of Hadoop

### Query Engines

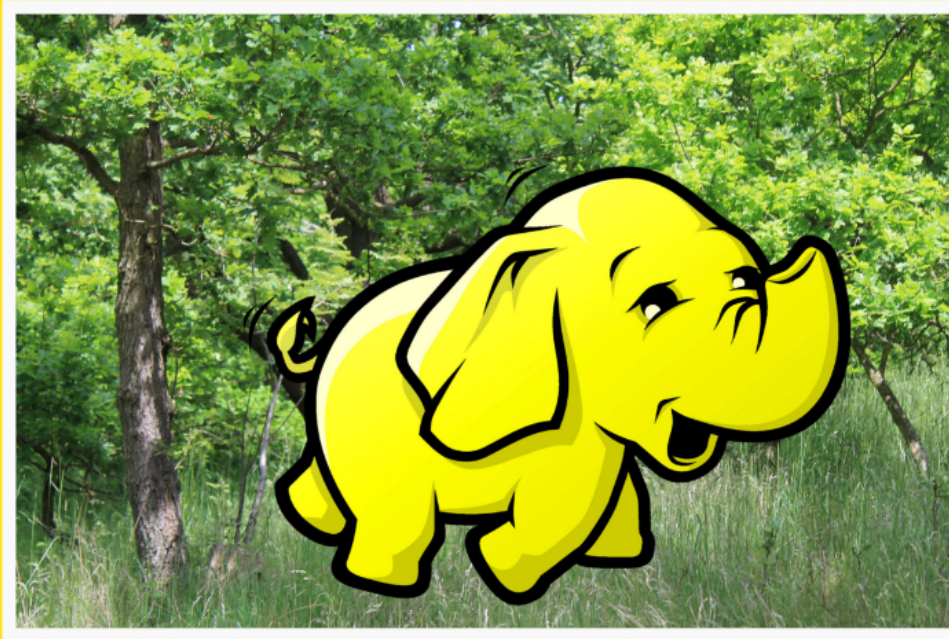### Core Hadoop Ecosystem

### External Data Storage

# Hadoop Overview

## What is Hadoop?

"an open source **software platform** for **distributed storage** and **distributed processing** of **very large data sets** on **computer clusters** built from commodity hardware" - *Hortonworks*

## Why Hadoop?

- Data's too darn big - terabytes per day
- Vertical scaling doesn't cut it
  - Disk seek times
  - Hardware failures
  - Processing times
- Horizontal scaling is linear
- Hadoop: It's not just for batch processing anymore

## Hadoop History

- Google published GFS and MapReduce papers in 2003-2004
- Yahoo! was building "Nutch," an open source web search engine at the same time
- Hadoop was primarily driven by Doug Cutting and Tom White in 2006
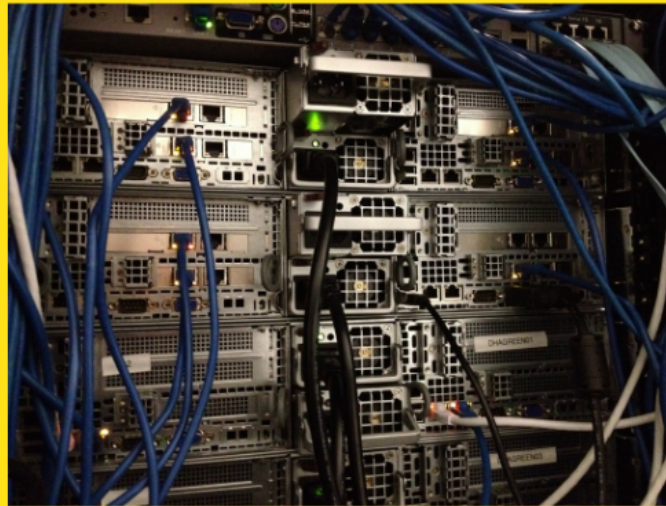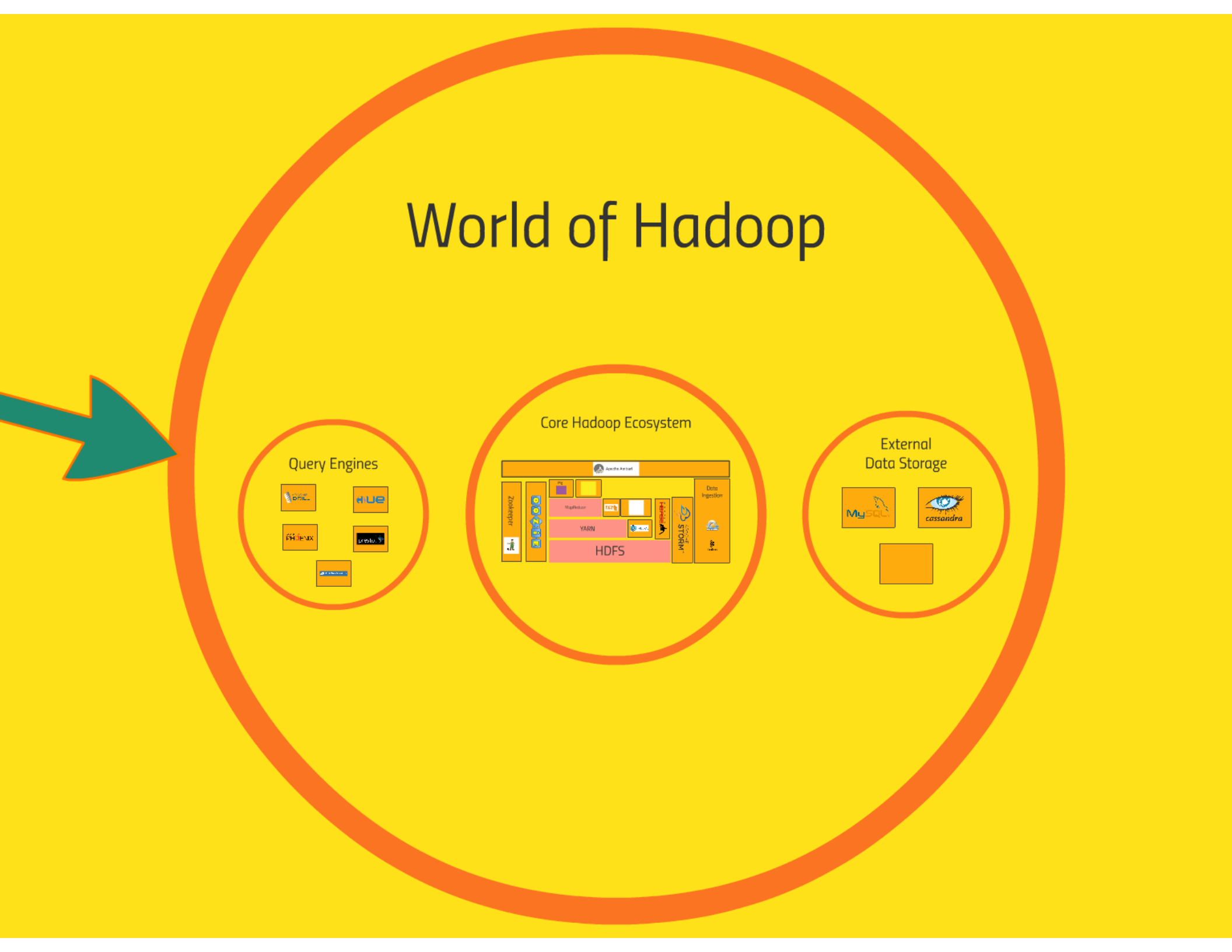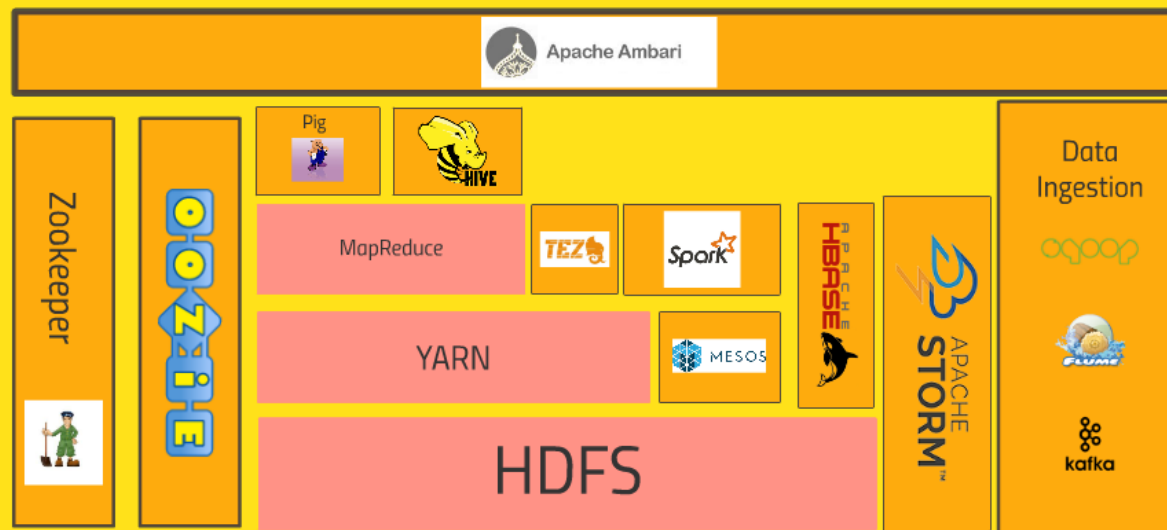- It's been evolving ever since...

# What is Hadoop?

"an open source **software platform** for **distributed storage** and **distributed processing** of **very large data sets** on **computer clusters** built from commodity hardware" - *Hortonworks*

# Hadoop History



- Google published GFS and MapReduce papers in 2003-2004
- Yahoo! was building "Nutch," an open source web search engine at the same time
- Hadoop was primarily driven by Doug Cutting and Tom White in 2006
- It's been evolving ever since...

# Why Hadoop?



- Data's too darn big - terabytes per day
- Vertical scaling doesn't cut it
  - Disk seek times
  - Hardware failures
  - Processing times
- Horizontal scaling is linear
- Hadoop: It's not just for batch processing anymore

World of Hadoop

Core Hadoop Ecosystem

MapReduce

TEZ

Spark

YARN

MESOS

APACHE HBASE

HDFS

MapReduce

YARN

# Pig

# Core Hadoop Ecosystem

Apache Ambari

Zookeeper

OOZIE

Pig

HIVE

MapReduce

TEZ

Spark

YARN

MESOS

HDFS

Apache HBASE

APACHE STORM

Data Ingestion

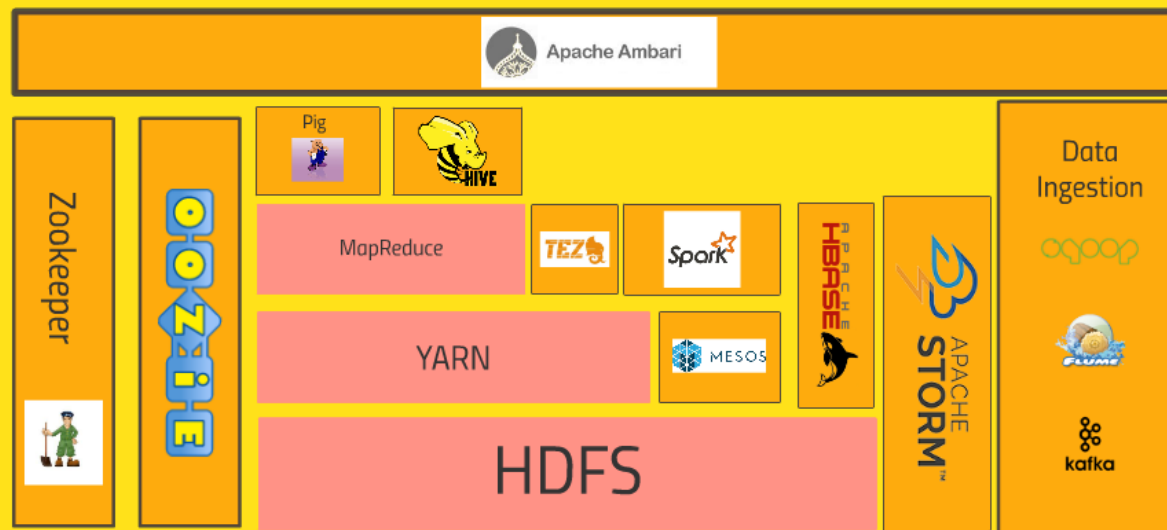sqoop

FLUME

kafka

# Core Hadoop Ecosystem

Apache Ambari

Core Hadoop Ecosystem

Core Hadoop Ecosystem

# Core Hadoop Ecosystem
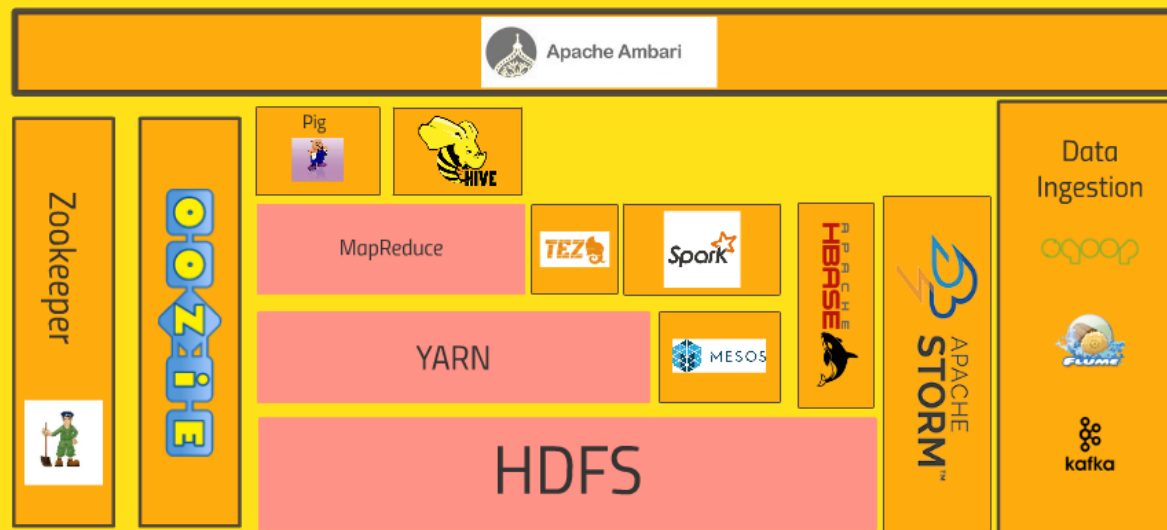
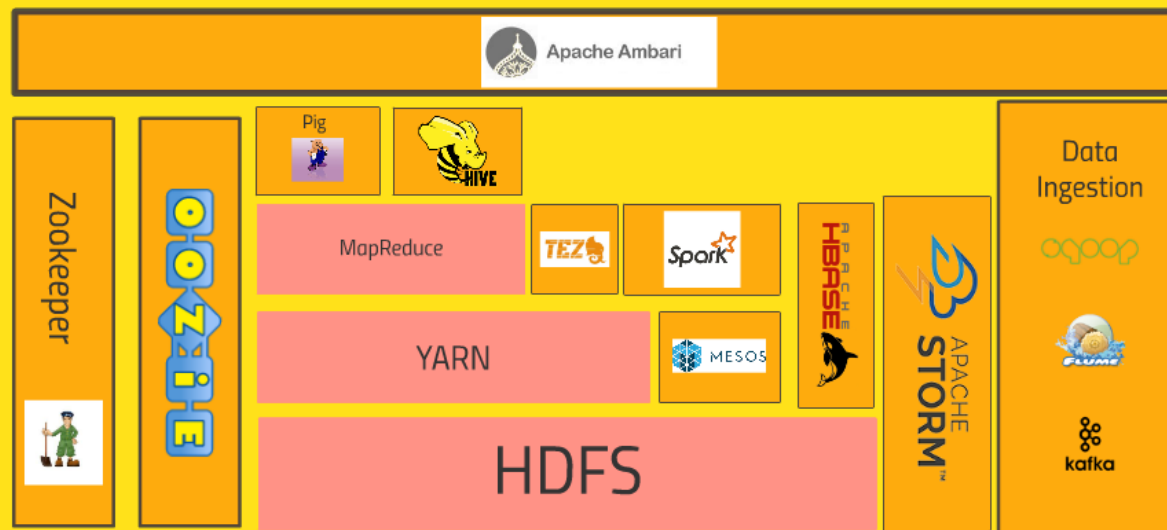# Core Hadoop Ecosystem

# Core Hadoop Ecosystem

# APACHE STORM™

# Core Hadoop Ecosystem

Zookeeper

OOZIE

Pig

MapReduce

YAR

Pig

Zookeeper

OOZIE

M

# Data Ingestion

APACHE HBASE

APACHE STORM™

SQOOP

FLUME

kafka

# Core Hadoop Ecosystem

# External Data Storage

cassandra

# External
# Data Storage

# Query Engines

APACHE DRILL

# World of Hadoop

## Core Hadoop Ecosystem

### Query Engines

### External Data Storage

# Core Hadoop Ecosystem

# External
# Data Storage

# Query Engines