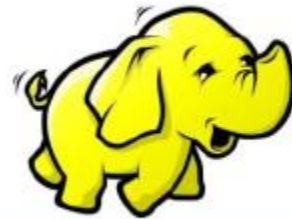# Agenda

✓ **Big Data Use Cases**
  ➢ Big Data Sources
  ➢ Common Big Data Customers Scenarios
  ➢ Hidden Treasure
  ➢ What Big Companies Have to Say..
  ➢ Hadoop Users in Detail

✓ **Motivation for Hadoop**
  ➢ What problems exists with traditional large scale system
  ➢ What requirement an alternative approach should have
  ➢ How Hadoop Addresses those requirements
  ➢ A brief Hadoop History
  ➢ Hadoop Core components
  ➢ Hadoop Key Characteristics

✓ **Hadoop Ecosystem**
  ➢ The components that creates Hadoop Eco-System

# Big Data Source

✓ Lots of Data(Terabytes or Petabytes)

✓ Systems / Enterprises generate huge amount of data from Terabytes to and even Petabytes of information.



**A airline jet collects 10 terabytes of sensor data for every 30 minutes of flying time.**
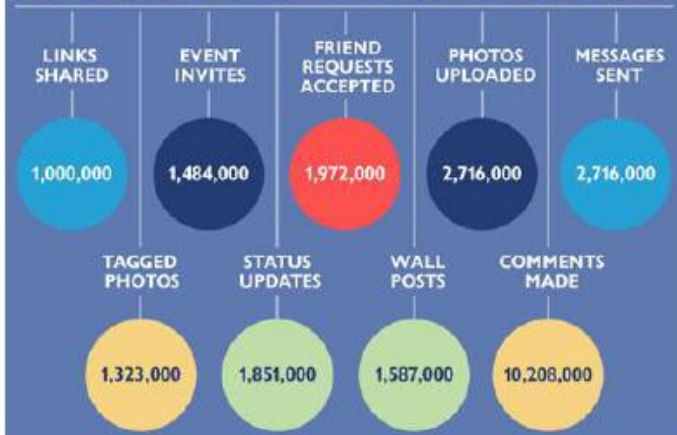


**NYSE generates about one terabyte of new trade data per day to Perform stock trading analytics to determine trends for optimal trades.**

# Facebook Example

# Twitter Example



✓ Twitter has over **500 million** registered users.

✓ The USA, whose **141.8 million** accounts represents 27.4 percent of all Twitter users, good enough to finish well ahead of Brazil, Japan, the UK and Indonesia.

✓ **79%** of US Twitter users are more like to recommend brands they follow .

✓ **67%** of US Twitter users are more likely to buy from brands they follow .

✓ **57%** of all companies that use social media for business use Twitter.

# What is Big Data?



✓ **IBM's definition – Big Data Characteristics**

http://www-01.ibm.com/software/data/bigdata/

Characteristics Of Big Data

Volume — 12 Terabytes of Tweets created each day

Velocity — Scrutinizes 5 million trade events created each day to identify potential fraud

Variety — Sensor data, audio, video, click streams, log files and more

# Un-Structured Data Is Exploding



- 2,500 exabytes of new information in 2012 with Internet as primary driver
- Digital universe grew by 62% last year to 800K petabytes and will grow to 1.2 "zettabytes" this year

# Common Big Data Customer Scenarios

✓ Web and e-tailing
  ✓ Recommendation Engines
  ✓ Ad Targeting
  ✓ Search Quality
  ✓ Abuse and Click Fraud Detection

✓ Telecommunications
  ✓ Customer Churn Prevention
  ✓ Network Performance Optimization
  ✓ Calling Data Record (CDR) Analysis
  ✓ Analyzing Network to Predict Failure

http://wiki.apache.org/hadoop/PoweredBy

# Common Big Data Customer Scenarios(Contd.)

✓ Government
  ✓ Fraud Detection And Cyber Security
  ✓ Welfare schemes
  ✓ Justice

✓ Healthcare & Life Sciences
  ✓ Health information exchange
  ✓ Gene sequencing
  ✓ Serialization
  ✓ Healthcare service quality improvements
  ✓ Drug Safety

http://wiki.apache.org/hadoop/PoweredBy

# Common Big Data Customer Scenarios(Contd.)

✓ **Banks and Financial services**
- ✓ Modeling True Risk
- ✓ Threat Analysis
- ✓ Fraud Detection
- ✓ Trade Surveillance
- ✓ Credit Scoring And Analysis

JPMorganChase

# Hidden Treasure

Case Study: Sears Holding Corporation
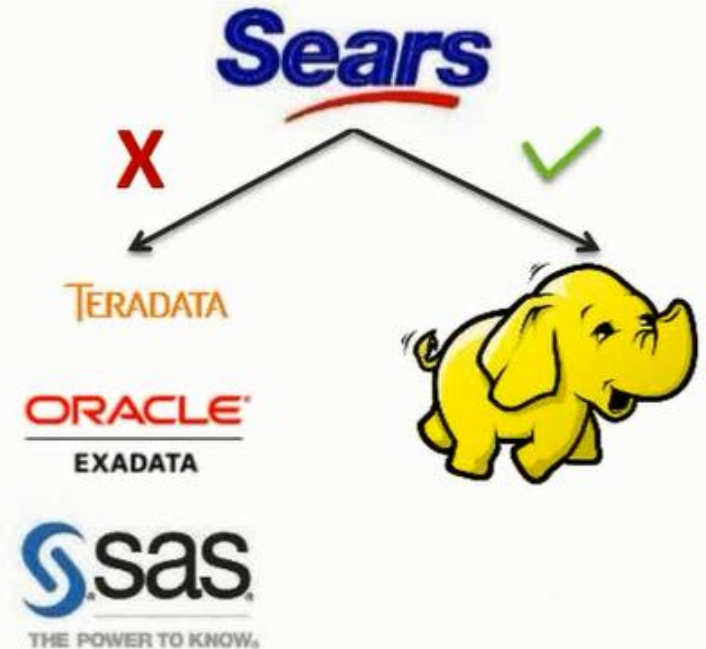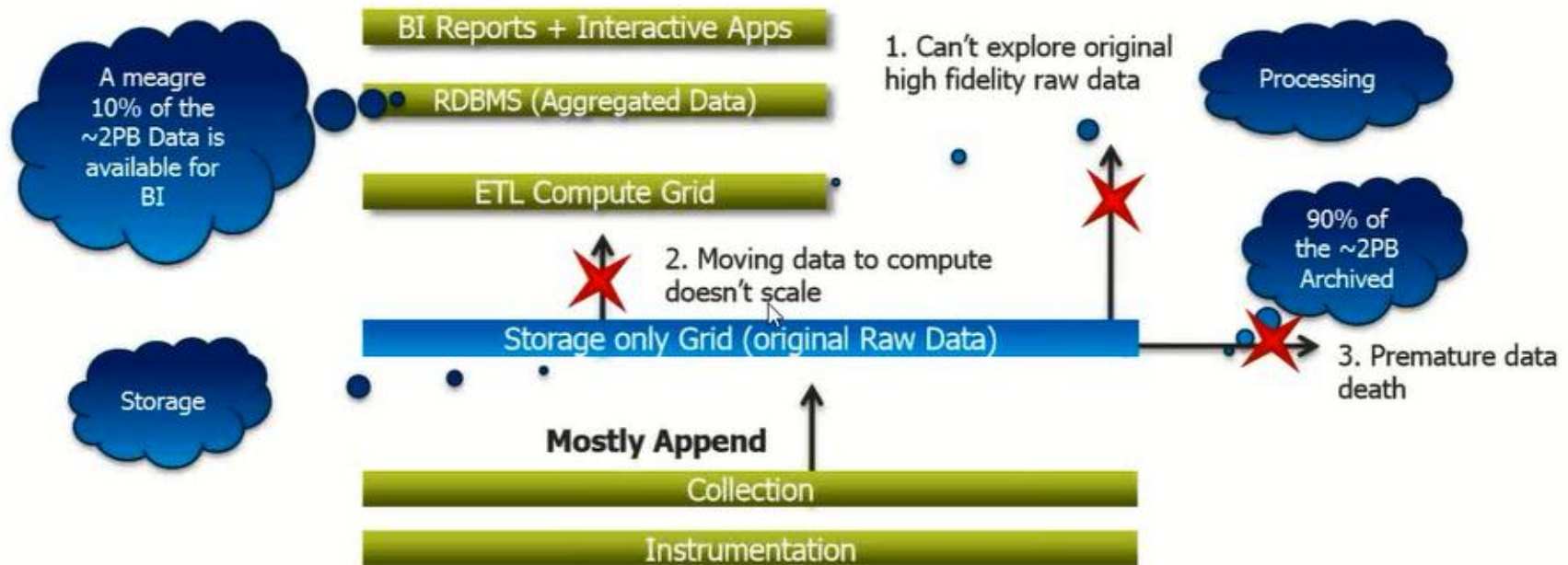
✓ Insight into data can provide **Business Advantage**.

✓ Some key early indicators can mean **Fortunes to Business**.

✓ **More Precise Analysis** with more data.

*Sears was using traditional systems such as Oracle Exadata, Teradata and SAS etc. to store and process the customer activity and sales data.

# Limitations of Existing Data Analytics Architecture



BI Reports + Interactive Apps

A meagre 10% of the ~2PB Data is available for BI

RDBMS (Aggregated Data)

1. Can't explore original high fidelity raw data

Processing

ETL Compute Grid

2. Moving data to compute doesn't scale

Storage only Grid (original Raw Data)

90% of the ~2PB Archived

Storage

3. Premature data death

Mostly Append

Collection

Instrumentation

http://www.informationweek.com/it-leadership/why-sears-is-going-all-in-on-hadoop/d/d-id/1107038?

Sears

# Solution: Using Hadoop



BI Reports + Interactive Apps

1. Data Exploration & Advanced analytics

RDBMS (Aggregated Data)

No Data Archiving

Entire ~2PB Data is available for processing

2. Scalable throughput for ETL & aggregation

Hadoop : Storage + Compute Grid

3. Keep data alive forever

Both Storage And Processing

Mostly Append

Collection

Instrumentation

*Sears moved to a 300-Node Hadoop cluster to keep 100% of its data available for processing rather than a meagre 10% as was the case with existing Non-Hadoop solutions.

Sears

# What Big Companies Have To Say…

**McKinsey**

*"Analyzing Big Data sets will become a key basis for competition."*

*"Leaders in every sector will have to grapple the implications of Big Data."*

**Gartner**

*"Big Data analytics are rapidly emerging as the preferred solution to business and technology trends that are disrupting."*

*"Enterprises should not delay implementation of Big Data Analytics."*

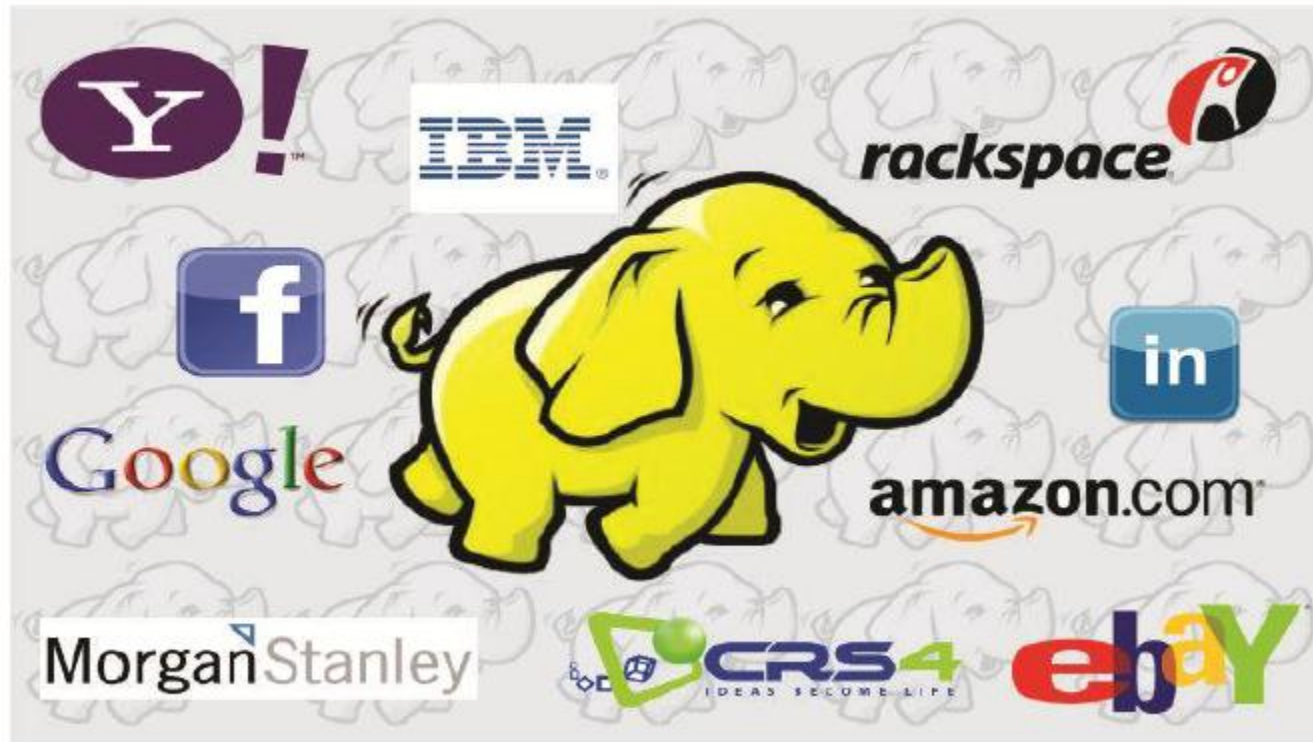*"Use Hadoop to gain a competitive advantage over more risk-averse enterprises."*

**Forrester Research**

*"Prioritize Big Data projects that might benefit from Hadoop."*

# Some of the Hadoop Users

# Hadoop Users – In Detail

**AOL**
- We use Hadoop for variety of things ranging from ETL style processing and statistics generation to running advanced algorithms for doing behavioral analysis and targeting.
- The Cluster that we use for mainly behavioral analysis and targeting has 150 machines, Intel Xeon, dual processors, dual core, each with 16GB Ram and 800 GB hard-disk.

**EBay**
- 532 nodes cluster (8 * 532 cores, 5.3PB).
- Heavy usage of Java MapReduce, Pig, Hive, HBase
- Using it for Search optimization and Research.

**Facebook**
- We use Hadoop to store copies of internal log and dimension data sources and use it as a source for reporting/analytics and machine learning.
- Currently we have 2 major clusters:
  - A 1100-machine cluster with 8800 cores and about 12 PB raw storage.
  - A 300-machine cluster with 2400 cores and about 3 PB raw storage.
  - Each (commodity) node has 8 cores and 12 TB of storage.
  - We are heavy users of both streaming as well as the Java APIs. We have built a higher level data warehousing framework using these features called Hive (see the http://hadoop.apache.org/hive/). We have also developed a FUSE implementation over HDFS.

# Hadoop Users – In Detail

**LinkedIn**
- *We have multiple grids divided up based upon purpose.*
- Hardware:
  - ~800 Westmere-based HP SL 170x, with 2x4 cores, 24GB RAM, 6x2TB SATA
  - ~1900 Westmere-based SuperMicro X8DTT-H, with 2x6 cores, 24GB RAM, 6x2TB SATA
  - ~1400 Sandy Bridge-based SuperMicro with 2x6 cores, 32GB RAM, 6x2TB SATA

**Openstat**
- Hadoop is used to run a customizable web analytics log analysis and reporting
- 50-node production workflow cluster (dual quad-core Xeons, 16GB of RAM, 4-6 HDDs) and a couple of smaller clusters for individual analytics purposes
- About 500 mln of events processed daily, 15 bln monthly
- Cluster generates about 25 GB of reports daily

**Rackspace**
- 30 node cluster (Dual-Core, 4-8GB RAM, 1.5TB/node storage)

**Telenav**
- 60-Node cluster for our Location-Based Content Processing including machine learning algorithms for Statistical Categorization, Deduping, Aggregation & Curation (Hardware: 2.5 GHz Quad-core Xeon, 4GB RAM, 13TB HDFS storage).
- Private cloud for rapid server-farm setup for stage and test environments.(Using Elastic N-Node cluster)
- Public cloud for exploratory projects that require rapid servers for scalability and computing surges (Using Elastic N-Node cluster)
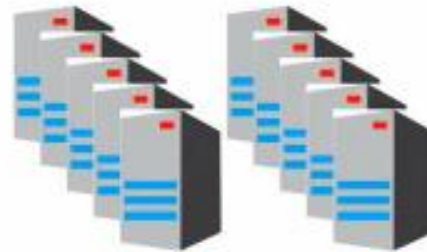
# Why DFS?



Read 1 TB Data

**1 Machine**
- 4 I/O Channels
- Each Channel – 100 MB/s

**10 Machines**
- 4 I/O Channels
- Each Channel – 100 MB/s

# Why DFS?

# Why DFS?
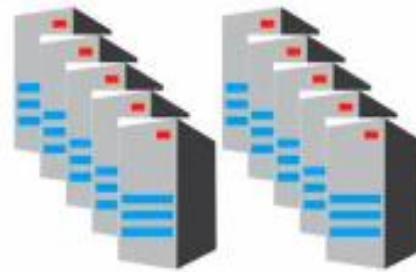
**Read 1 TB Data**

| 1 Machine | 10 Machines |
|---|---|
| • 4 I/O Channels<br>• Each Channel – 100 MB/s | • 4 I/O Channels<br>• Each Channel – 100 MB/s |
| 45 Minutes | 4.5 Minutes |

# What is DFS?



Before DFS consolidation

\\Chicago_server\Homedirs
\\Chicago_maxi\Projects
\\Houston_server\Reports
\\Denver_server\Software

After DFS consolidation

\\maxi-pedia.com\Public\Homedirs
\\maxi-pedia.com\Public\Projects
\\maxi-pedia.com\Public\Reports
\\maxi-pedia.com\Public\Software

# Motivation For Hadoop (Contd.)

✓ Traditional Distributed Systems - Problems

➤ Programming for traditional distributed system is complex

- Data exchange requires synchronization
- Finite bandwidth is available
- It is difficult to deal with partial failures of the system

➤ Ken Arnold,CORBA designer

- "Failure is the defining difference between distributed and local programming, So you have to design a distributed systems with the expectation of failure"
- Developers spend more time designing for failure than they do actually working on the problem itself.

# Motivation For Hadoop (Contd.)

✓ Distributed Systems: Data Storage

➢ Typically , data for a distributed system is stored on SAN

➢ At Compute time, data is copied to the compute nodes

➢ Fine for relatively limited amounts of data

# Motivation For Hadoop (Contd.)

✓ Data becomes the Bottleneck

➢ Processing power doubles every two years

➢ Processing speed is no longer a problem

➢ Getting data to the processors becomes the bottleneck

➢ Quick Calculation

- Typical disk data transfer rate:75MB/sec

- Time Taken to transfer 100GB of data to the processor : approx. 22 minutes!

  ❑ Assuming sustained reads

  ❑ Actual time will be worse, since most servers have less than 100GB of RAM available,

- Significant amount of complex processing performed on that data

➢ A new approach is needed

# Motivation For Hadoop (Contd.)

✓ Requirements for a new approach(CAP principle)

➢ Partial Failure Support

- System must support partial failure
- Failure of the component should result in a graceful degradation of application performance. Not complete failure of the entire system.

➢ Data Recoverability and high availability

- If the component of the system fails, its workload should be assumed by still-functioning units in the system.
- Failure should not result in the loss of any data

➢ Component Recovery

- If the component of the system fails and then recovers ,it should be able to rejoin the system without requiring  a full restart of the system.

# Motivation For Hadoop (Contd.)

✓ Requirements for a new approach(CAP principle)

➢ Consistency
- Component Failures during the execution of a job should not affect the outcome of the job, data should be consistent across all the replicas
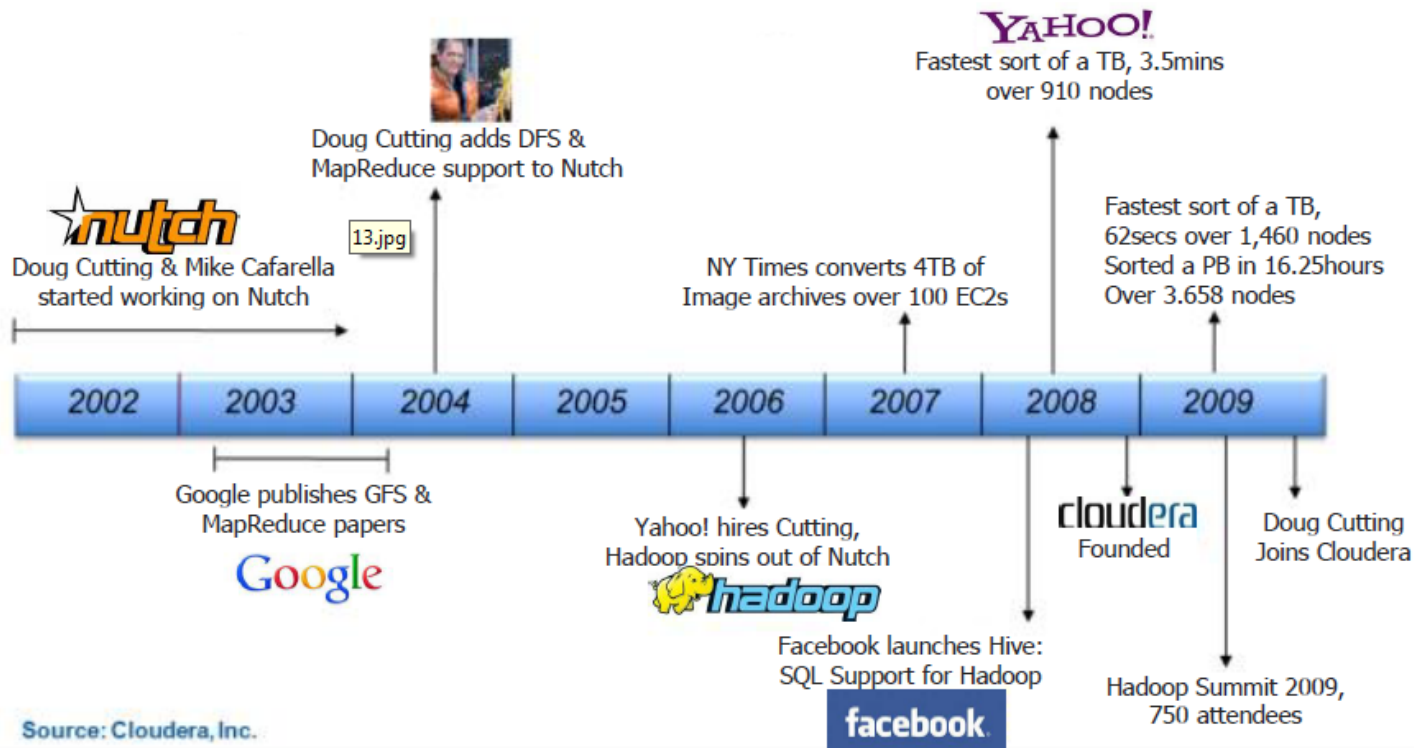
➢ Scalability
- Adding load to the system should result in a graceful decline in performance of individual jobs ,not failure of the system
- Increasing resources should support a proportional increase in load capacity

➢ Partition Tolerance
- Even if data is partitioned and store in multiple nodes (may be located across different geographic locations) , once you do hadoop fs -cat filename, it should display all the data from all the partitions seamlessly. Abstract this storage from user so that he will see as if data comes from the same machine.

# Introducing Hadoop

✓ Hadoop's History

# Introducing Hadoop

✓ Hadoop Core concept

➢ Apache Hadoop is a framework that allows for the **distributed processing** of large data sets **stored** across clusters of commodity computers using a **simple programming model**.

➢ It's an open source framework which supports **scaled out storage** and **distributed processing**

➢ Distribute the data as it is initially stored in the system
- Individual node can work on the data local to that node
- No data transfer over the network is required for initial processing

➢ Application is written in high level code
- Developers need not worry about network programming, temporal dependencies or low level infrastructures.

➢ Nodes talk to each other as little as possible
- Developers should not write any code which communicates between nodes – "Shared Nothing" architecture.

➢ Data is spread among machines in advance
- Computation happens where data is stored,wherevr possible
- Data is replicated multiple times on the system for increased availability and reliability
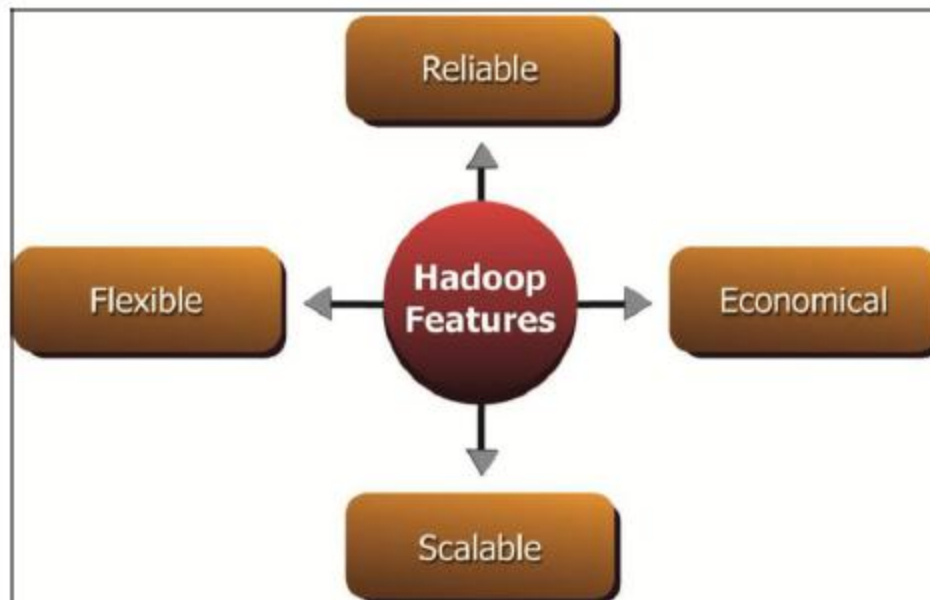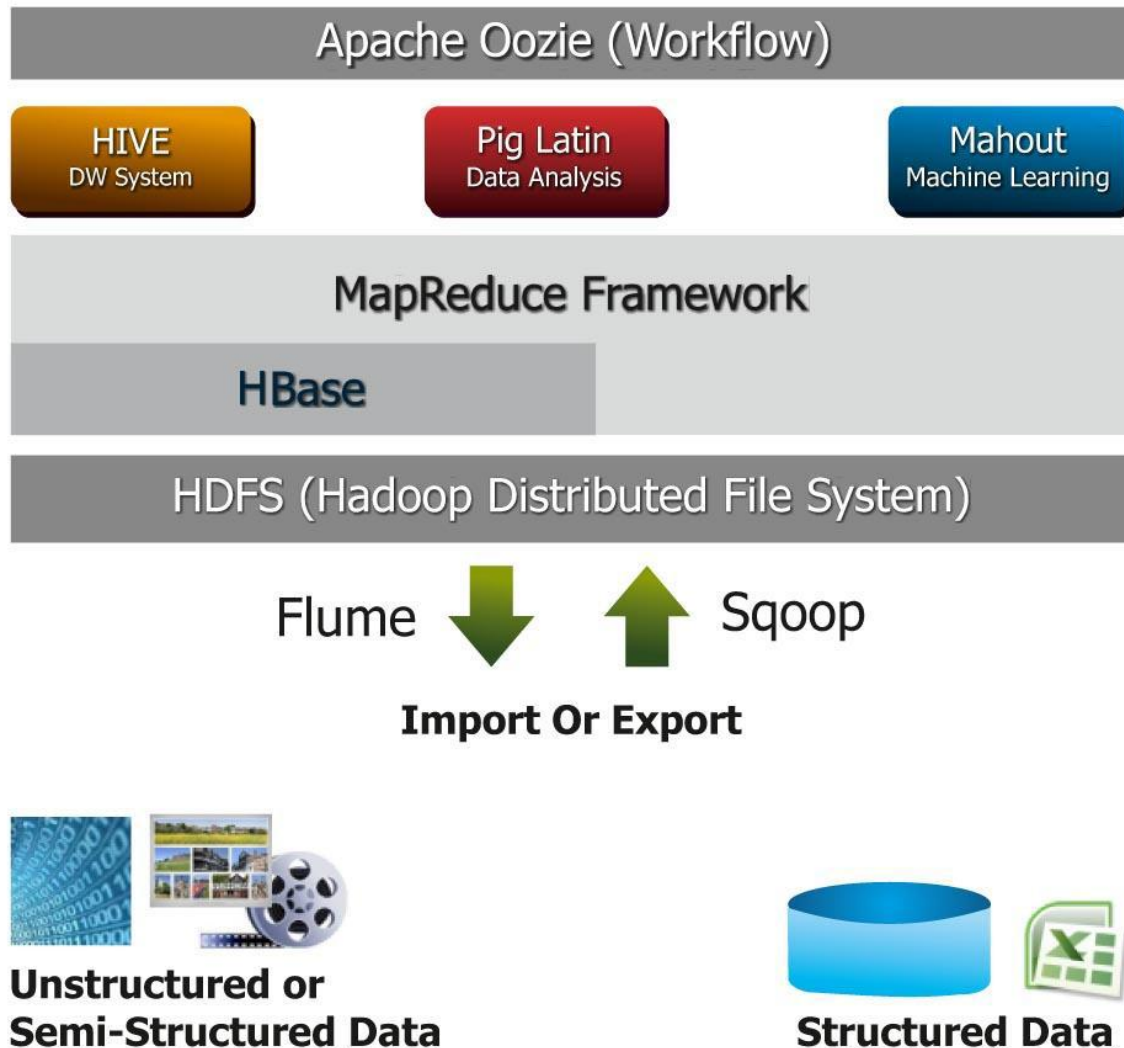
# Introducing Hadoop

✓ Hadoop Core Components

➢ Hadoop consists of two core components
- Hadoop Distributed File System (HDFS)(Storage)
- MapReduce(Processing)

➢ There are many other projects based around core Hadoop
- Often referred to as the 'Hadoop Ecosystem'
- Pig, Hive, HBase, Oozie, Sqoop, etc

➢ A set of machines running HDFS and MapReduce is known as a Hadoop Cluster
- Individual machines are known as nodes
- A cluster can have as few as one node, as many as several thousand
- More nodes=better performance!

# Hadoop Key Characteristics

# Hadoop Ecosystem

# Thank You