

What is Hadoop?



Agenda

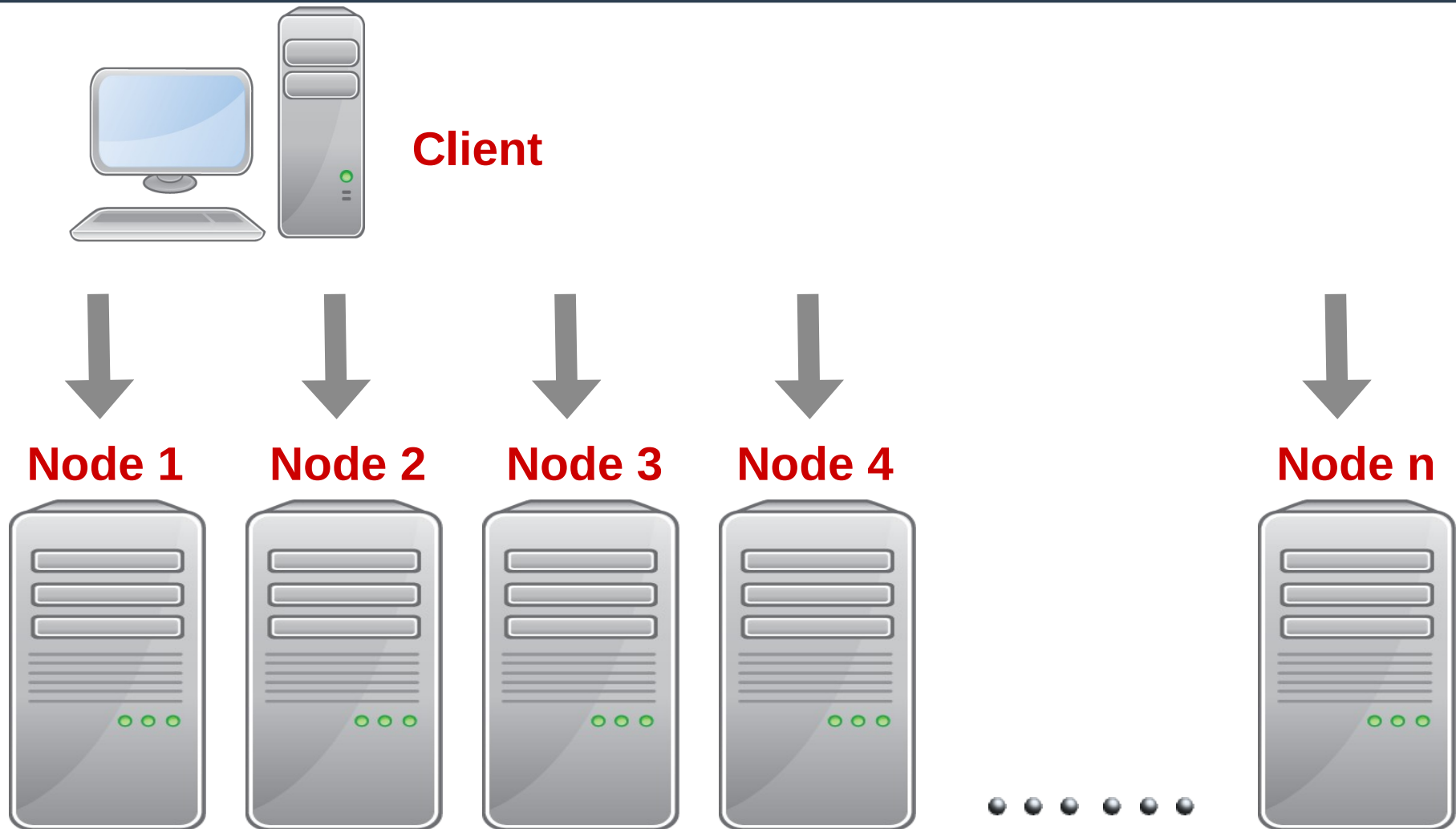
- **What is Hadoop?**
- **Hadoop Cluster**
- **History**
- **Hadoop System Principles**
- **Comparisons to RDBMS**
- **Hadoop Distribution Vendors**

What is Hadoop?



- Hadoop is a **reliable, scalable platform for storage and analysis**.
- It runs on commodity hardware.
- It is open source.

Hadoop Cluster



Hadoop Cluster

- **A set of "cheap" commodity hardware**
 - Servers Networked together
 - Performs same task as a system



Hadoop Cluster

- **No need for super-computers, It uses commodity hardware**
- **Not desktops**



History

- **Named after an elephant toy Started as a sub-project of Apache Nutch**
 - Nutch's job is to index the web and expose it for searching
 - Open Source alternative to Google
 - Started by Doug Cutting
- **In 2004 Google publishes Google File System(GFS) and MapReduce framework papers**
- **Doug Cutting and Nutch team implemented Google's frameworks in Nutch**
- **In 2006 Yahoo! hires Doug Cutting to work on Hadoop with a dedicated team**
 - In 2008 Hadoop became Apache Top Level Project

Hadoop System Principles

- **Scale-Out** rather than **Scale-Up**.
- Bring **code to data** rather than data to code.
- **Fault Tolerance**.
- **Abstract complexity** of distributed and concurrent applications.

Scale-Out Vs Scale-Up

- **Scale-up(vertical)**

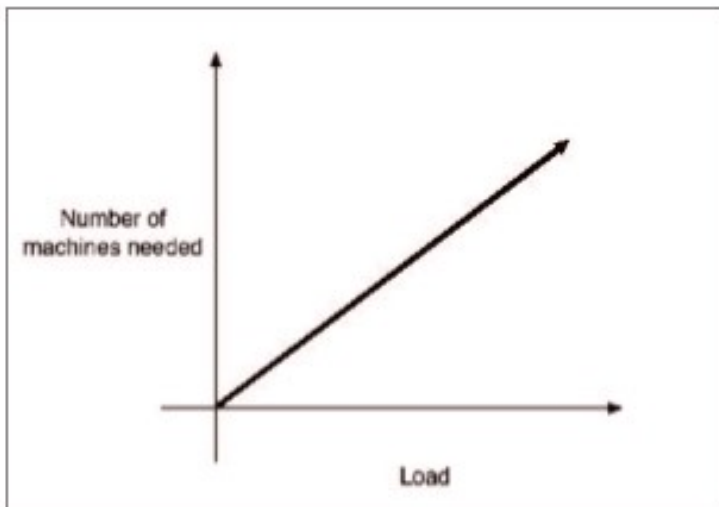
- Add additional resources to an existing node (CPU, RAM)
- Moore's Law can't keep up with data growth
- New units must be purchased if required resources can not be added
- Expensive and hard to implement

- **Scale-Out(Horizontal)**

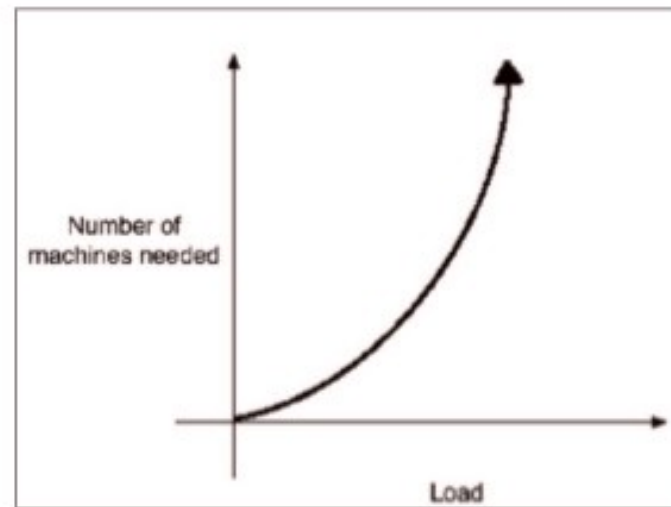
- Add more nodes/machines to an existing distributed application
- Software Layer is designed for node additions or removal
- Hadoop takes this approach - A set of nodes are bonded together as a single distributed system
- Very easy to scale down as well

Scalability

Linear vs. Non-Linear Scalability



Linear Scalability

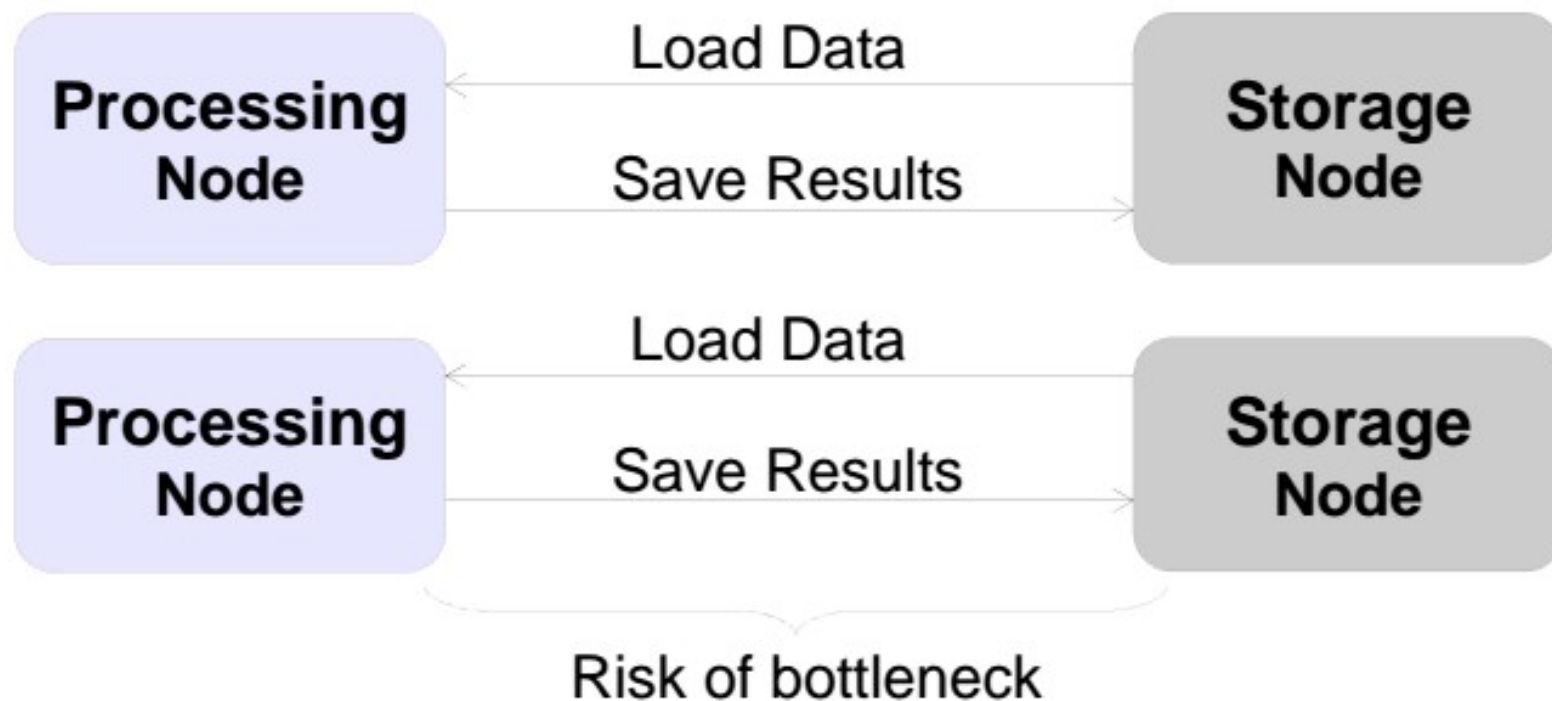


Non- Linear Scalability

"A linearly scalable system can maintain performance under increased load by adding resources in proportion to the increased load"

- **Traditional data processing architecture**
 - Nodes are divided into separate processing and storage, connected by high-capacity link
 - Many data-intensive applications are not CPU demanding causing bottlenecks in network

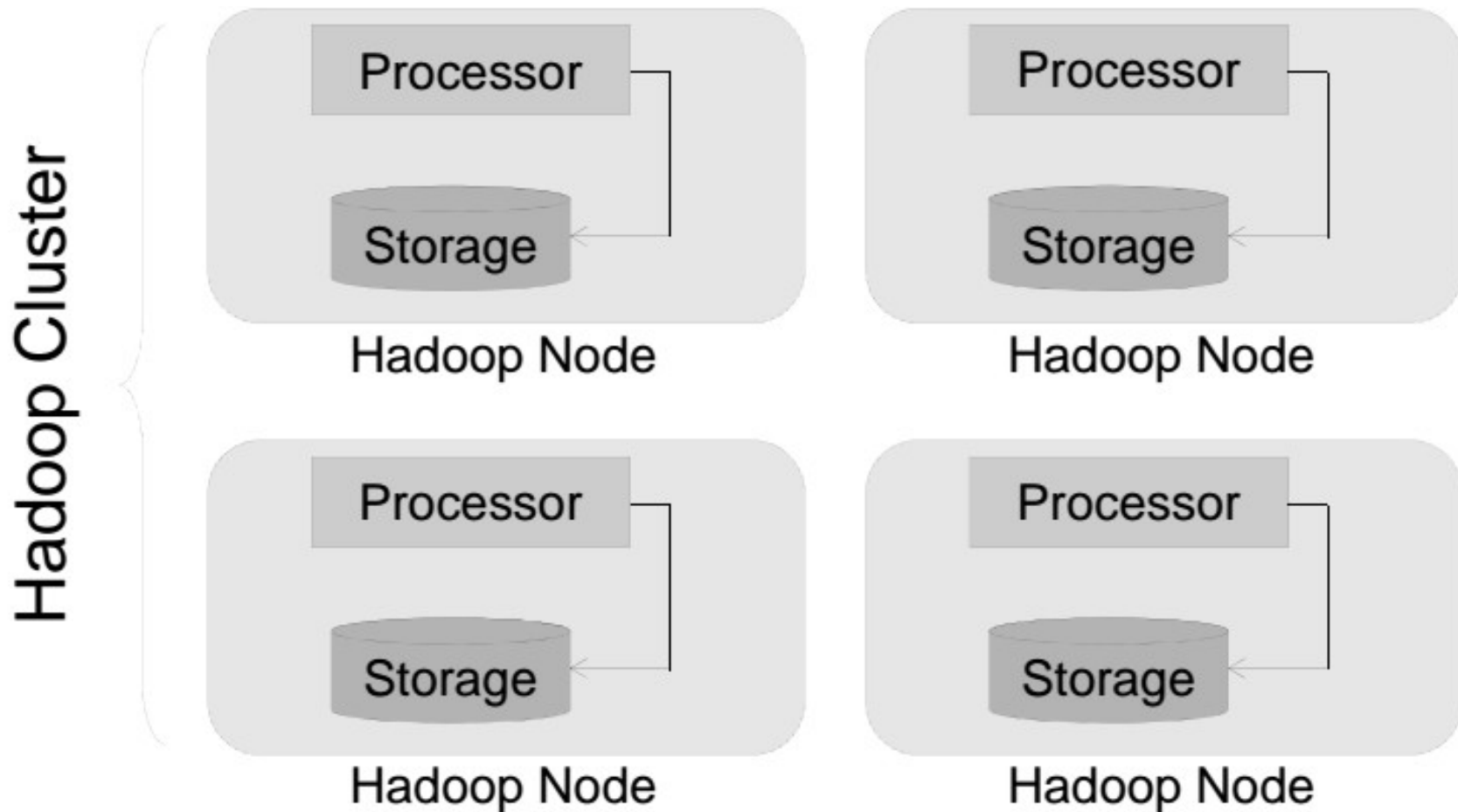
Data to Code



Code to Data

- **Hadoop co-locates processors and storage**
 - Code is moved to data (size is tiny, usually in Kbs)
 - Processors execute code and access underlying local storage

Code to Data



Fault Tolerance

- **Given a large number machines, failures are common**
- **Hadoop is designed to handle node failures**
 - Data is replicated
 - Tasks are retried

Abstract Complexity

- **Hadoop abstracts many complexities in distributed and concurrent applications**
 - **Defines small number of components**
 - **Provides simple and well defined interfaces of interactions between these components**
- **Frees developer from worrying about system level challenges**
 - **race conditions, data starvation**
 - **processing pipelines, data partitioning, code distribution etc.**
- **Allows developers to focus on application development and business logic**

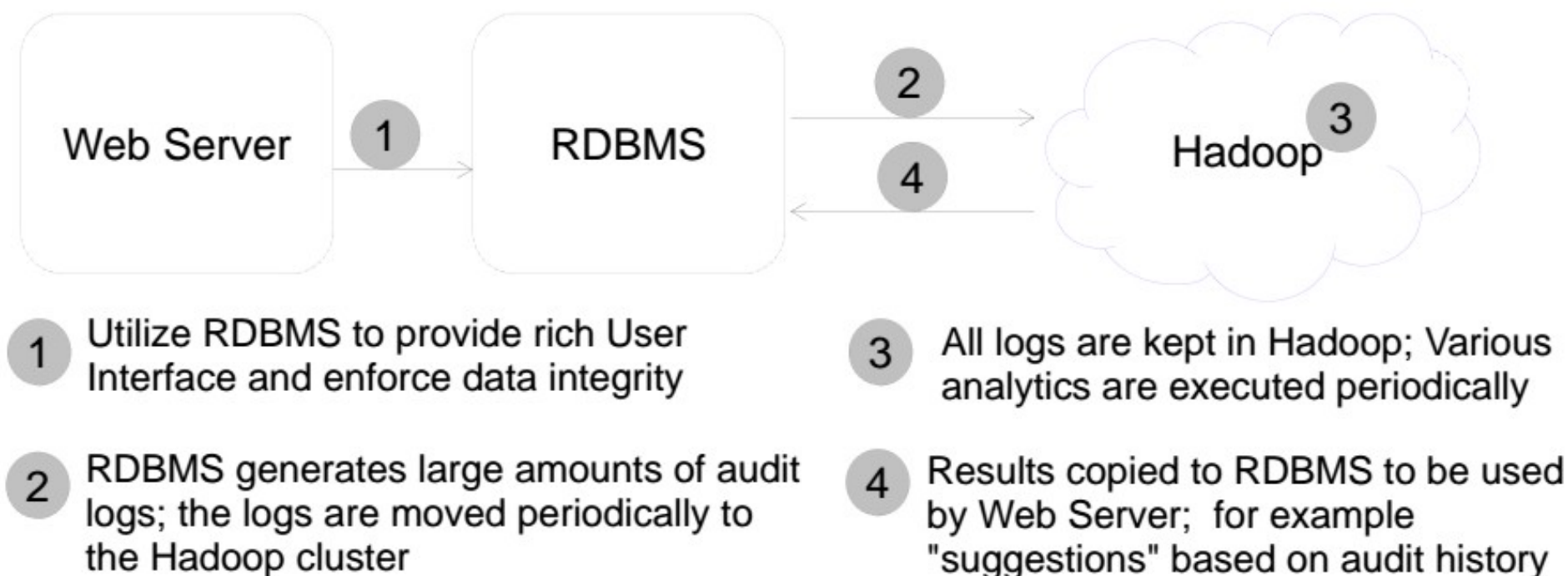
Comparisons to RDBMS

- **Hadoop and RDBMS frequently complement each other within an architecture**

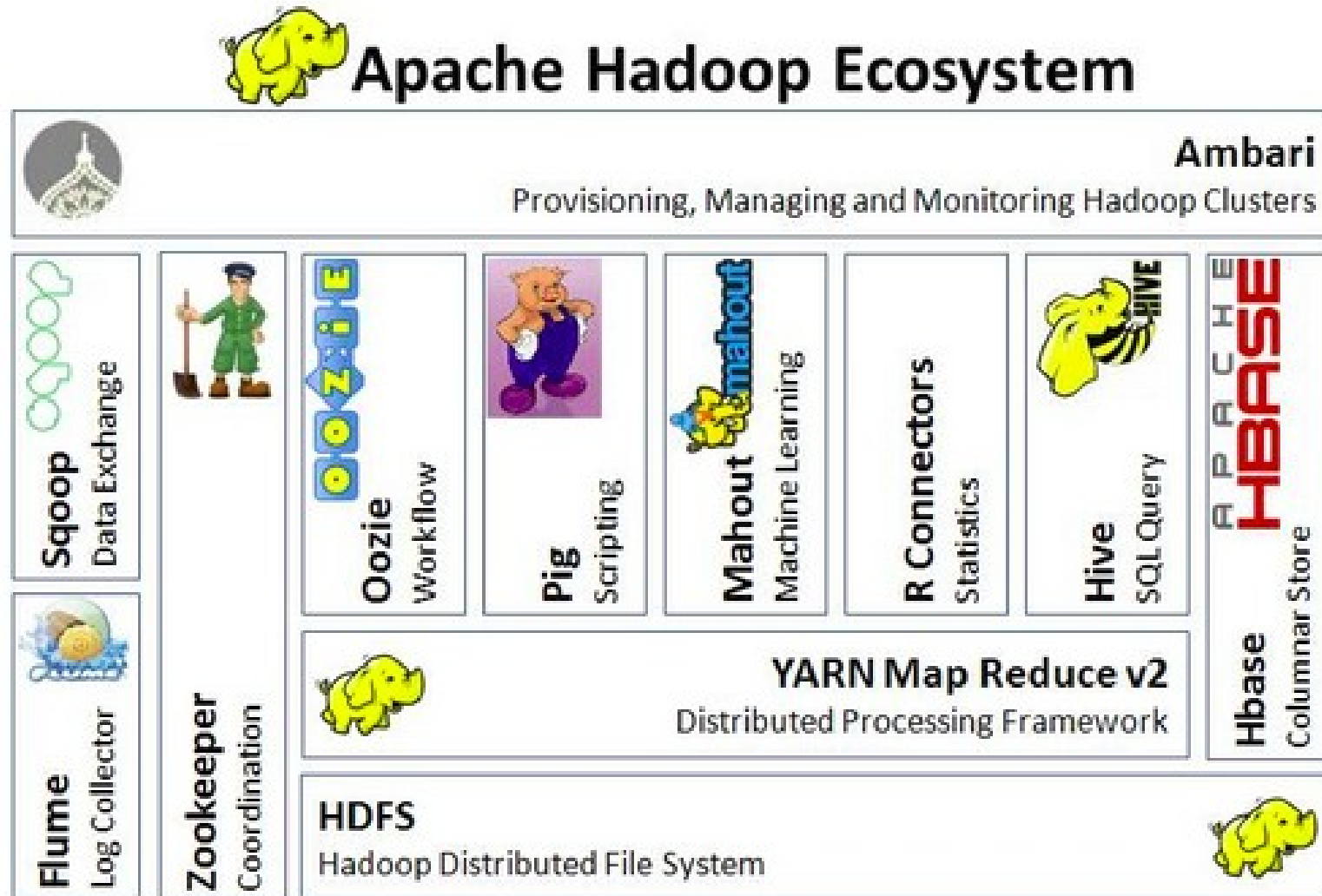
	Traditional RDBMS	MapReduce
Data size	Gigabytes	Petabytes
Access	Interactive and batch	Batch
Updates	Read and write many times	Write once, read many times
Transactions	ACID	None
Structure	Schema-on-write	Schema-on-read
Integrity	High	Low
Scaling	Nonlinear	Linear

Comparisons to RDBMS

- **For example, a website that**
 - has a small number of users
 - produces a large amount of audit logs



Hadoop Eco System



Hadoop Eco System

- **HDFS:** Hadoop Distributed FileSystem
- **MapReduce:** Distributed data processing framework
- **HBase:** Hadoop column database; supports batch and random reads and limited queries
- **Zookeeper:** Highly-Available Coordination Service
- **Oozie:** Hadoop workflow scheduler and manager
- **Pig:** Data processing language and execution environment
- **Hive:** Data warehouse with SQL interface

Hadoop Eco System

- **To start building an application, you need a file system**
 - In Hadoop world that would be **Hadoop Distributed File System (HDFS)**
- **Addition of a data store would provide a nicer interface to store and manage your data**
 - **HBase**: A key-value store implemented on top of HDFS
 - Traditionally one could use RDBMS on top of a local file system

Hadoop Eco System

HBase



Hadoop Distributed FileSystem (HDFS)

Hadoop Eco System

- **For batch processing, you will need to utilize a framework**
 - In Hadoop's world that would be **MapReduce**

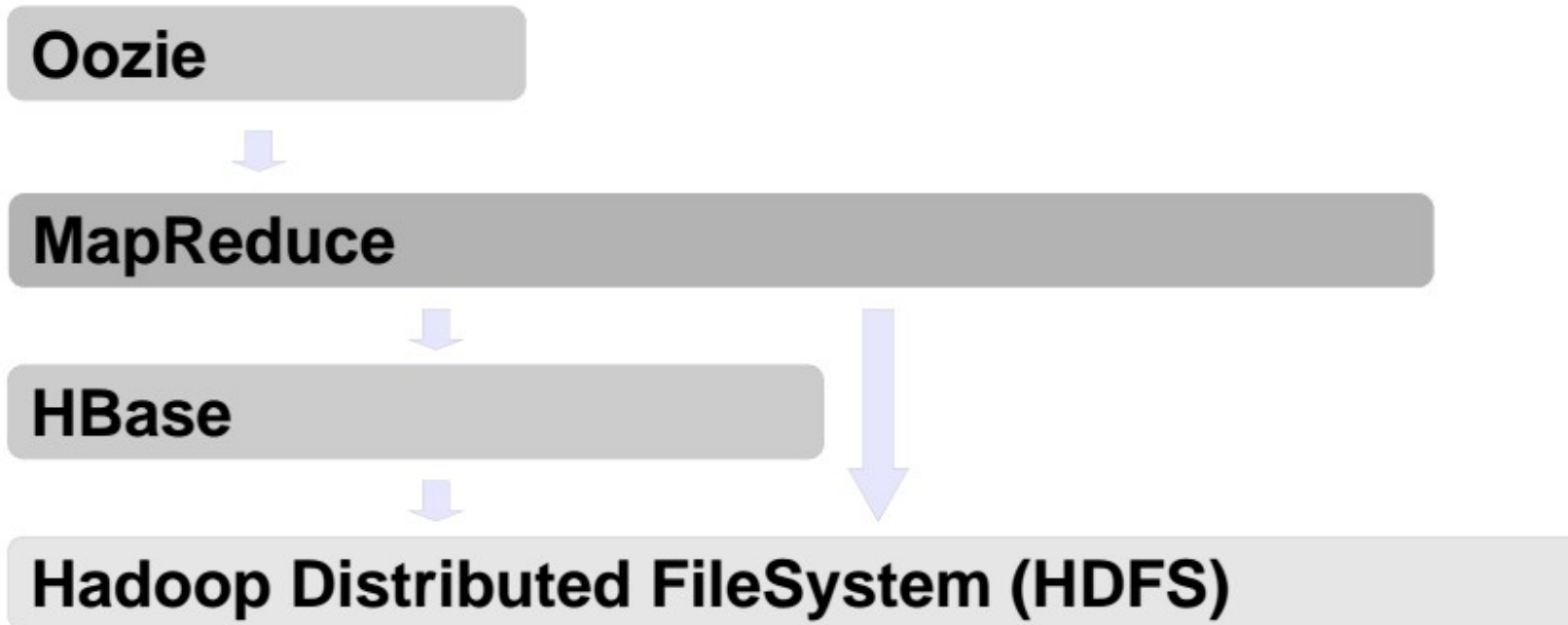
MapReduce

HBase

Hadoop Distributed FileSystem (HDFS)

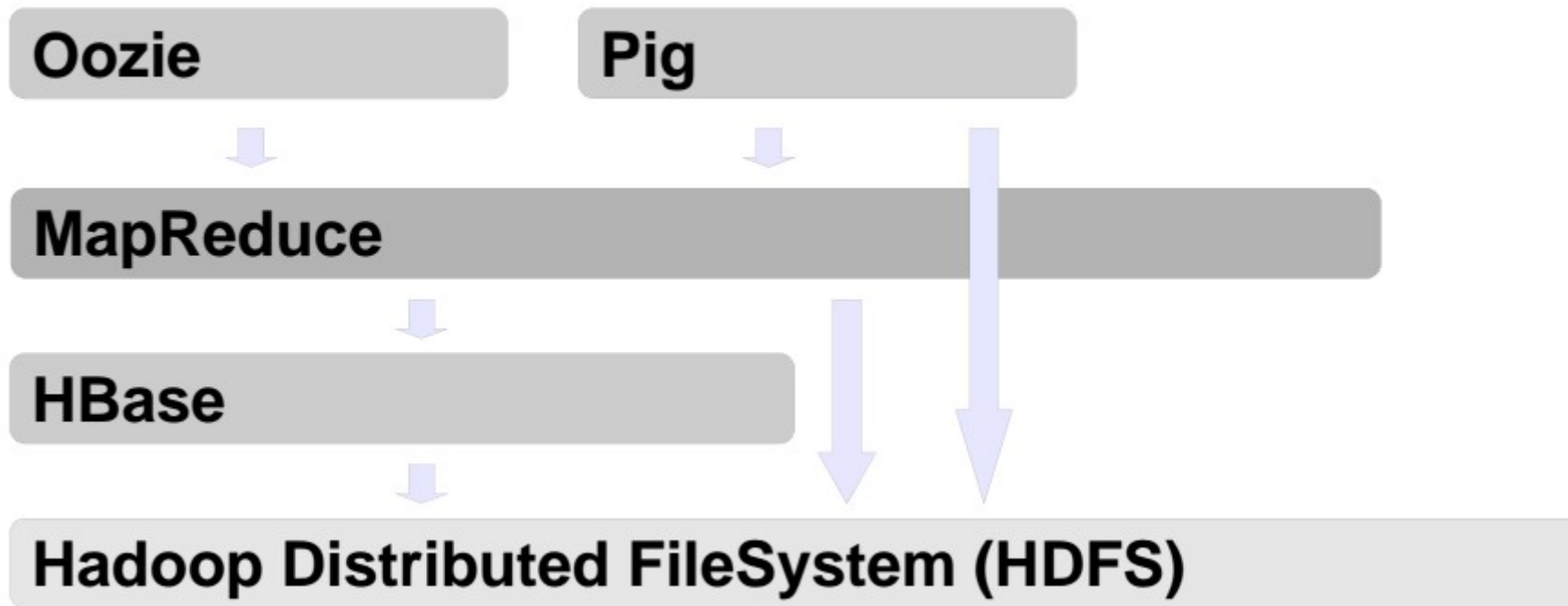
Hadoop Eco System

- **Many problems require MapReduce solution with multiple jobs**
 - **Apache Oozie** is a popular MapReduce workflow and coordination product



Hadoop Eco System

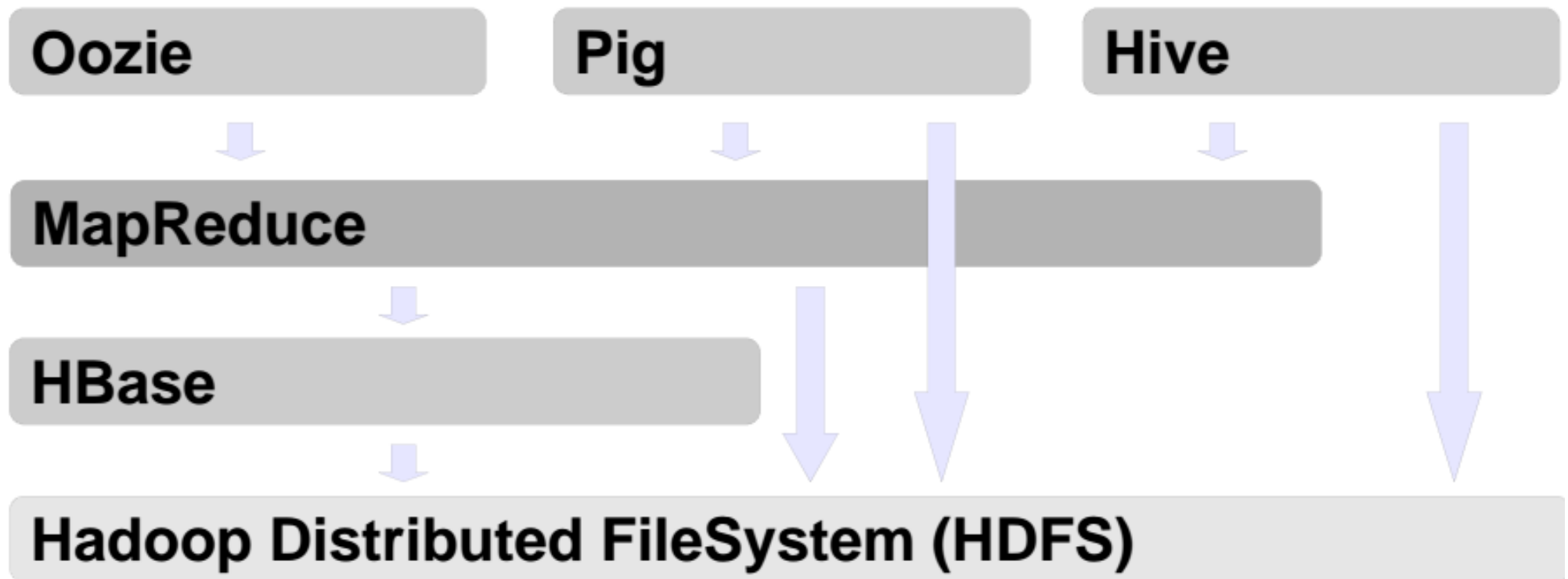
- **MapReduce paradigm may not work well for analysts and data scientists**
 - Addition of **Apache Pig**, a high-level data flow scripting language, may be beneficial



Hadoop Eco System

- **Your organization may have a good number of SQL experts**
 - Addition of **Apache Hive**, a data warehouse solution that provides a SQL based interface, may bridge the gap

Hadoop Eco System



Distribution Vendors

- **Cloudera** Distribution for Hadoop (CDH)
- **MapR** Distribution
- **Hortonworks** Data Platform(HDP)

cloudera®



Resources

- **Hadoop: The Definitive Guide**
 - Tom White (Author)
 - O'Reilly Media; 4th Edition.

