

# Introduction To Big Data

**<https://github.com/bobbylovemovie/trainbigdata>**

**bobbylovemovie@gmail.com**

**Big Data**  
**Big Data Analytics**  
**Data Science**  
**Machine Learning**  
**Artificial Intelligence**  
**Deep Learning**



# What is Big Data ?

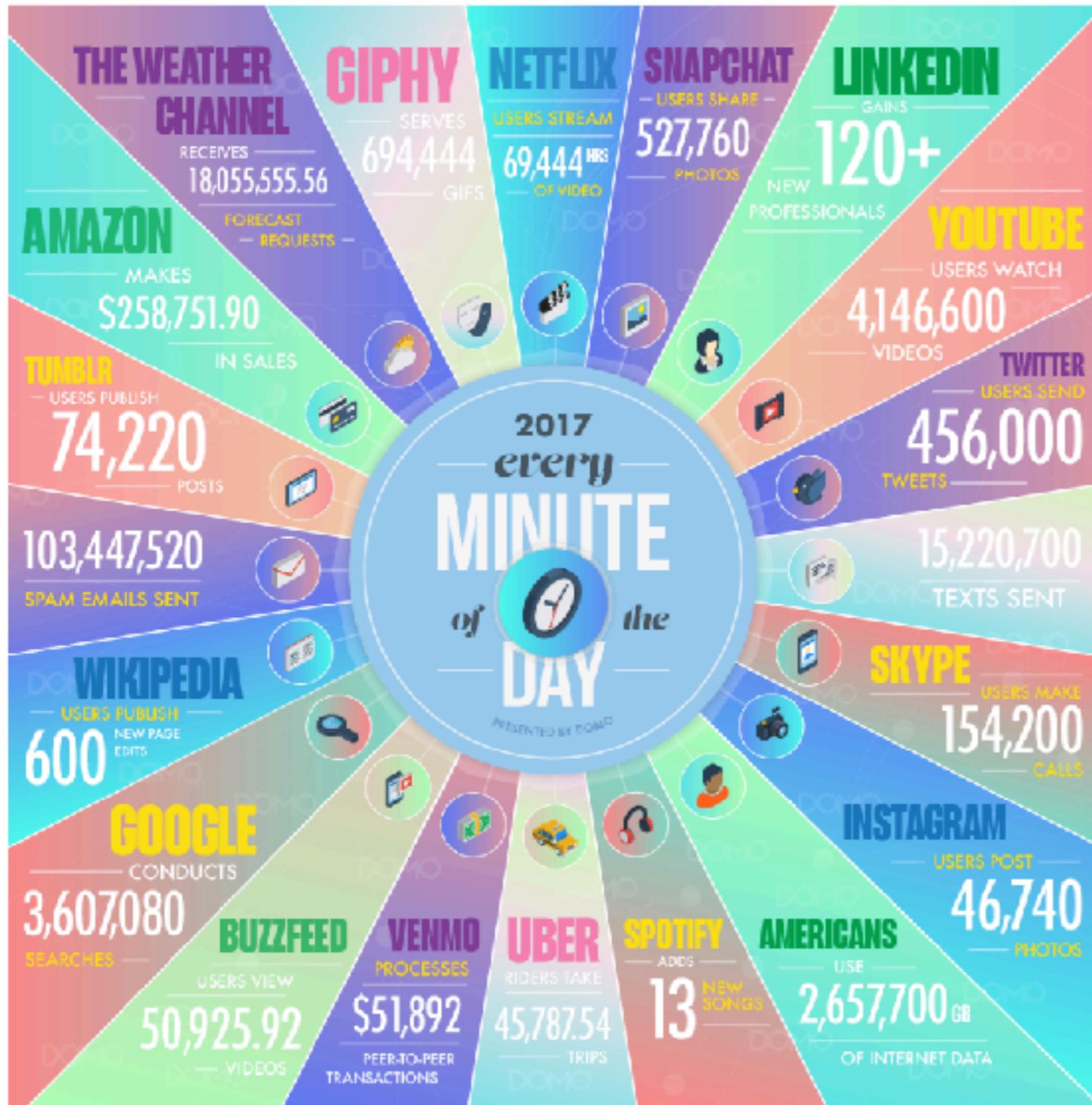
**“ Big data is data that exceeds the processing capacity of conventional database systems.**

**The data is too big, moves too fast, or doesn’t fit the structures of your database architectures.**

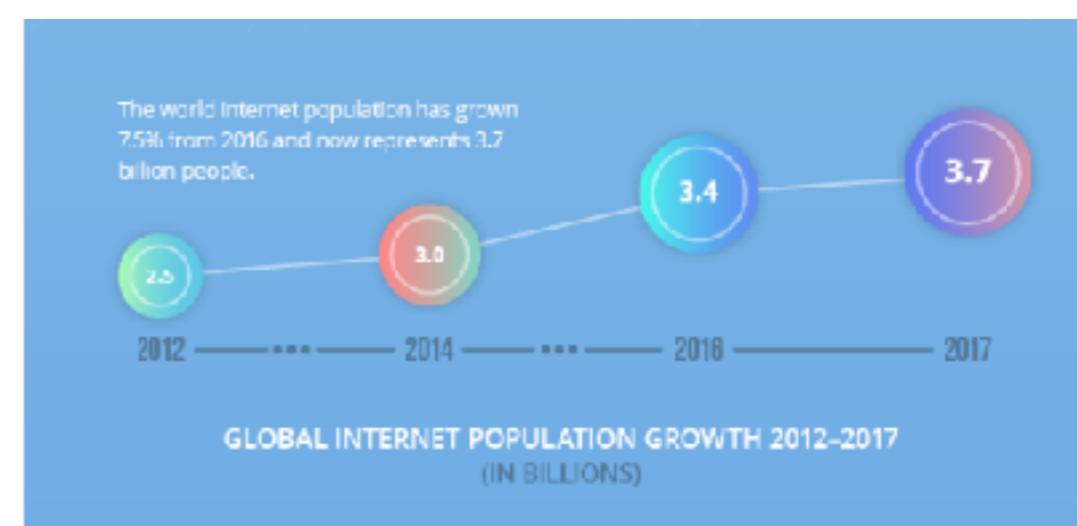
**To gain value from this data, you must choose an alternative way to process it. ”**

**Big Data Now: O'Reilly Media**

# Why Big Data ?



## DATA NEVER SLEEP 4.0



# Three Characteristics of Big Data

## Volume

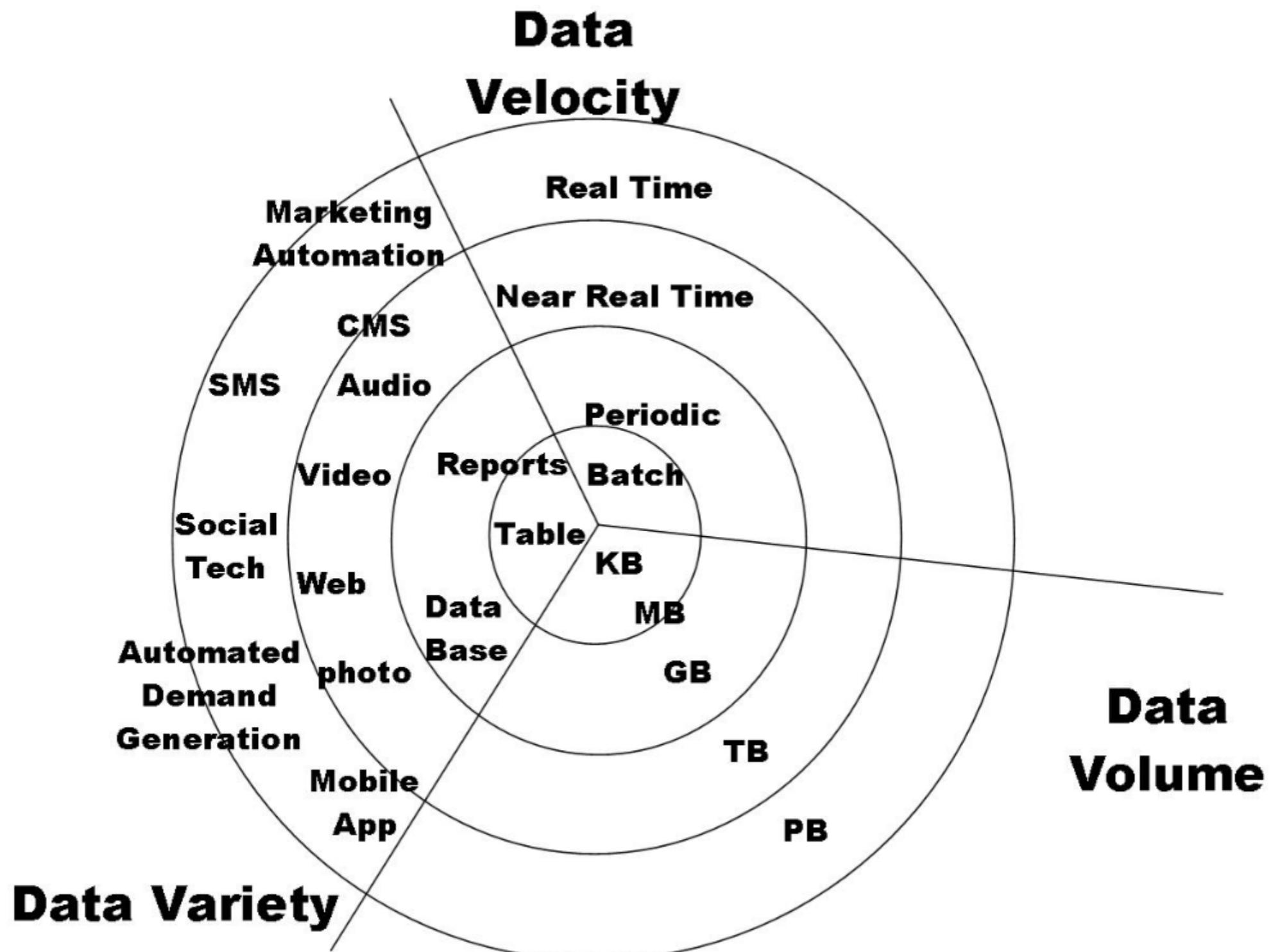
- Volumes of data are larger than those conventional relational database infrastructures can cope with

## Velocity

- Rate at which data flows in is much faster.
  - Mobile event and interaction by users.
  - Video, image , audio from users

## Variety

- the source data is diverse, and doesn't fall into neat relational structures eg. text from social networks, image data, a raw feed directly from a sensor source.



# Big Data = Volume, Variety and Velocity (3Vs)

Amount of new data stored varies across geography

New data stored<sup>1</sup> by geography, 2010  
Petabytes



<sup>1</sup> New data stored defined as the amount of available storage used in a given year; see appendix for more on the definition and assumptions.

SOURCE: IDC storage reports; McKinsey Global Institute analysis

## Velocity

30 billion pieces of content are shared on Facebook every month.



4 billion hours of video are watched on YouTube each month



400 Million Tweets are sent per day  
200M monthly active users



Source: IRM

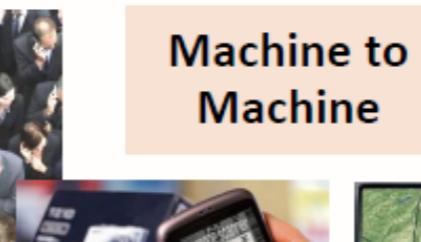
## Variety



People to Machine



People to People



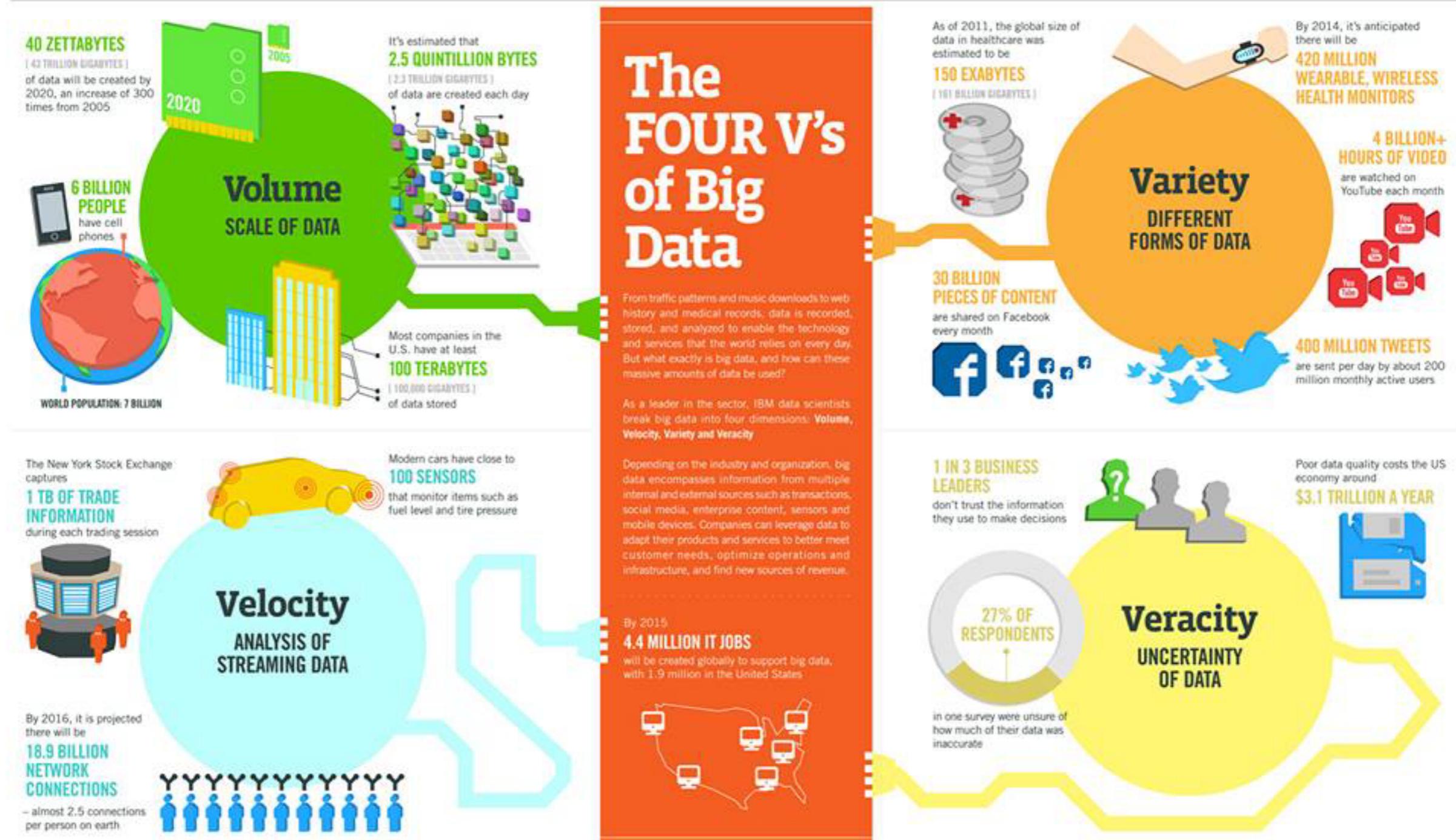
Machine to Machine



## Scale

1000 kilobytes = 1 Megabyte  
1000 Megabytes = 1 Gigabyte  
1000 Gigabytes = 1 Terabyte  
1000 Terabytes = 1 Petabyte  
1000 Petabytes = 1 Exabyte  
1000 Exabytes = 1 Zettabyte  
1000 Zettabytes = 1 Yottabyte  
1000 Yottabytes = 1 Bronobyte  
1000 Bronobytes = 1 Geopbyte

# 4Vs of Big Data



## 40 ZETTABYTES

[ 43 TRILLION GIGABYTES ]  
of data will be created by  
2020, an increase of 300  
times from 2005



## 6 BILLION PEOPLE

have cell  
phones



WORLD POPULATION: 7 BILLION



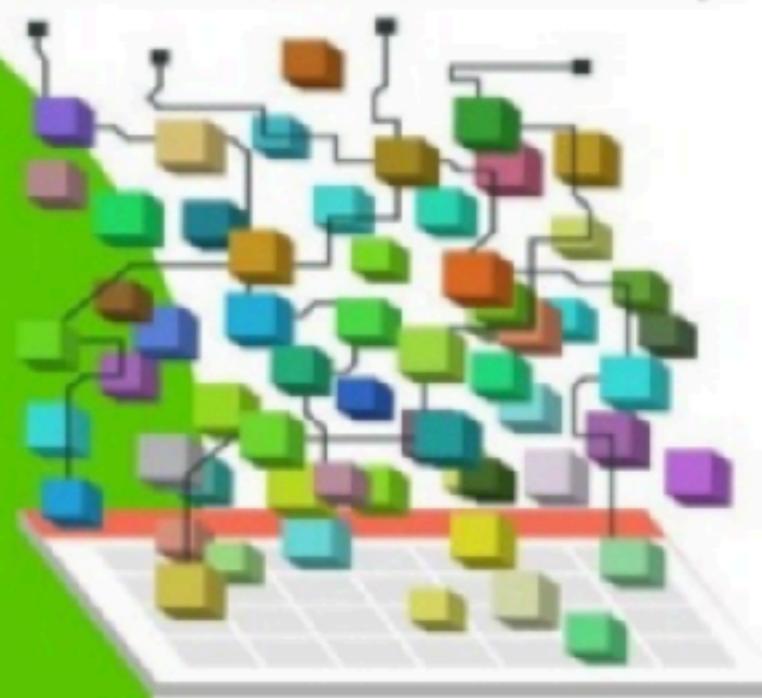
## Volume SCALE OF DATA



It's estimated that

## 2.5 QUINTILLION BYTES

[ 2.3 TRILLION GIGABYTES ]  
of data are created each day



Most companies in the  
U.S. have at least

## 100 TERABYTES

[ 100,000 GIGABYTES ]  
of data stored

The New York Stock Exchange captures

**1 TB OF TRADE INFORMATION**

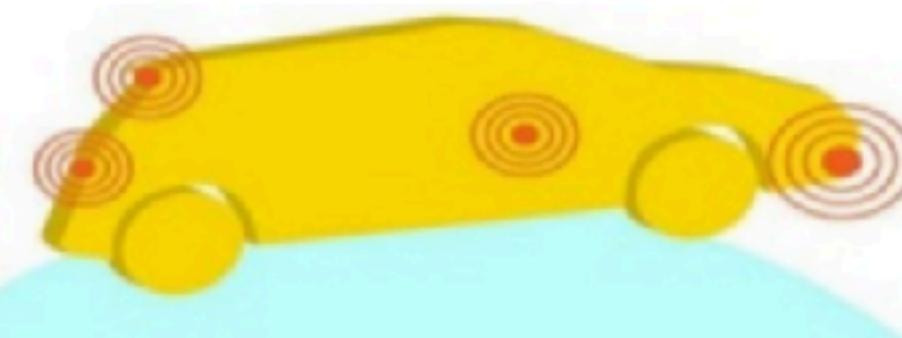
during each trading session



By 2016, it is projected there will be

**18.9 BILLION NETWORK CONNECTIONS**

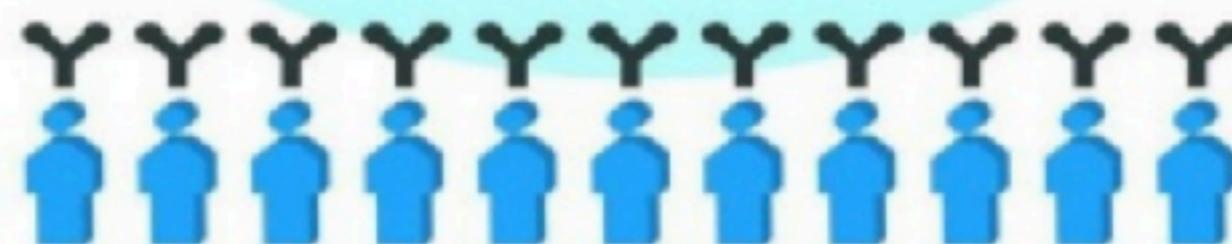
– almost 2.5 connections per person on earth



Modern cars have close to **100 SENSORS**

that monitor items such as fuel level and tire pressure

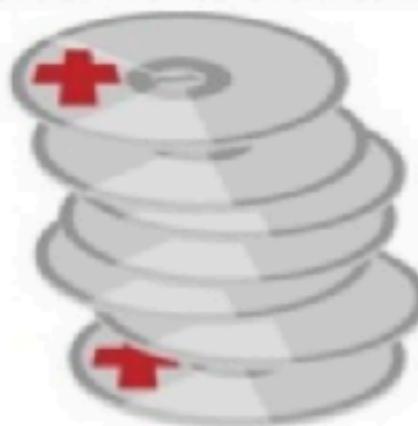
## Velocity ANALYSIS OF STREAMING DATA



As of 2011, the global size of data in healthcare was estimated to be

**150 EXABYTES**

[ 161 BILLION GIGABYTES ]



**30 BILLION  
PIECES OF CONTENT**

are shared on Facebook every month



## Variety DIFFERENT FORMS OF DATA

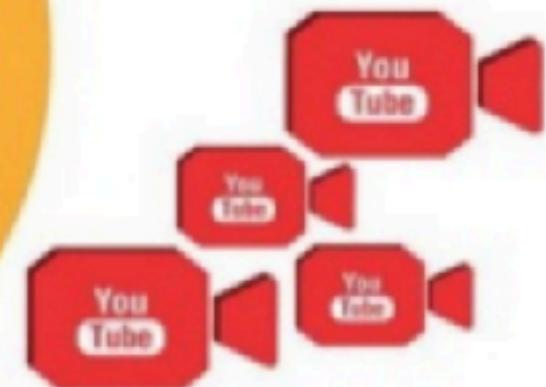


By 2014, it's anticipated there will be

**420 MILLION  
WEARABLE, WIRELESS  
HEALTH MONITORS**

**4 BILLION+  
HOURS OF VIDEO**

are watched on YouTube each month



**400 MILLION TWEETS**

are sent per day by about 200 million monthly active users

## 1 IN 3 BUSINESS LEADERS

don't trust the information  
they use to make decisions

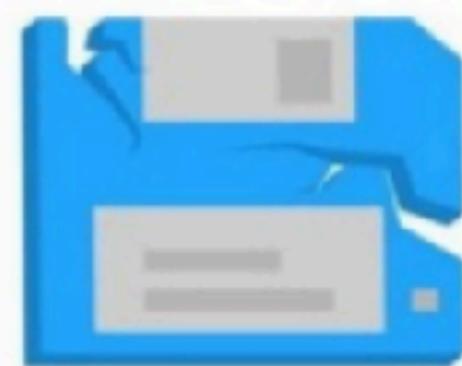


27% OF  
RESPONDENTS

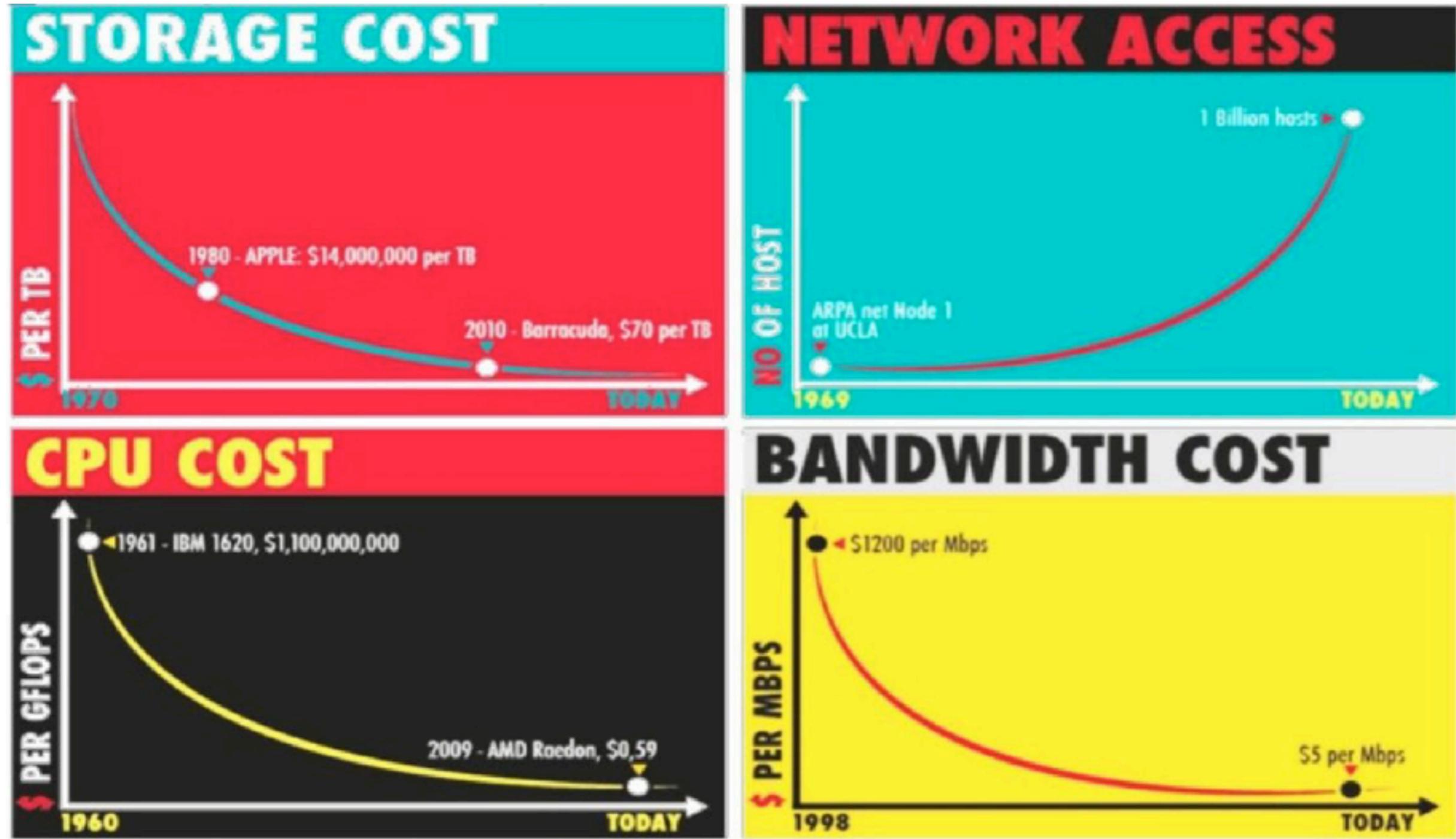
in one survey were unsure of  
how much of their data was  
inaccurate

# Veracity UNCERTAINTY OF DATA

Poor data quality costs the US  
economy around  
**\$3.1 TRILLION A YEAR**



# Big Data : Why Now?



Source: William EL KAIM, Enterprise Architecture and Technology Innovation 21

# DATA LAKE

# DATA LAKE



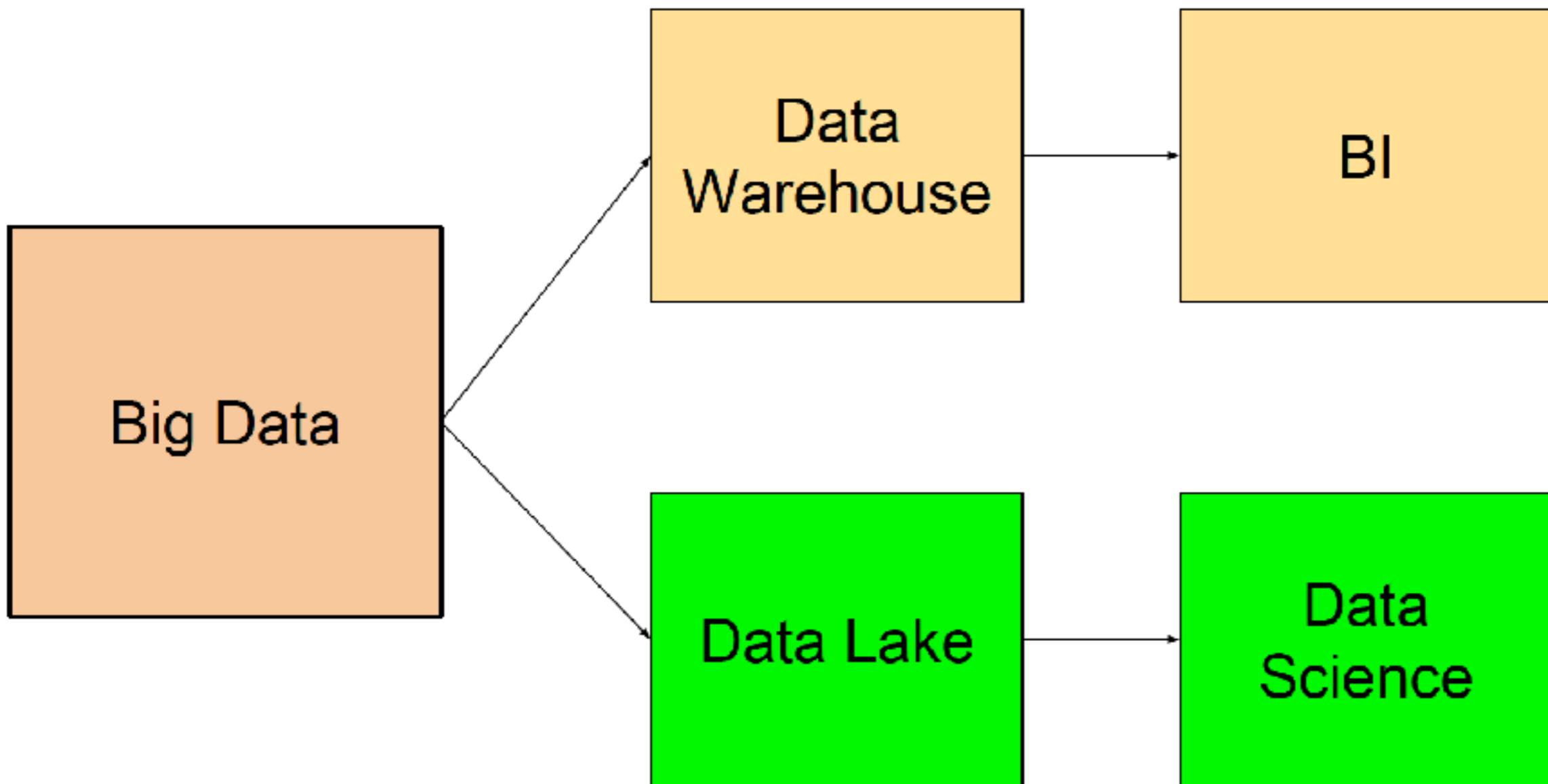
# The old way: Ask, then collect



# The new way: Collect, then ask

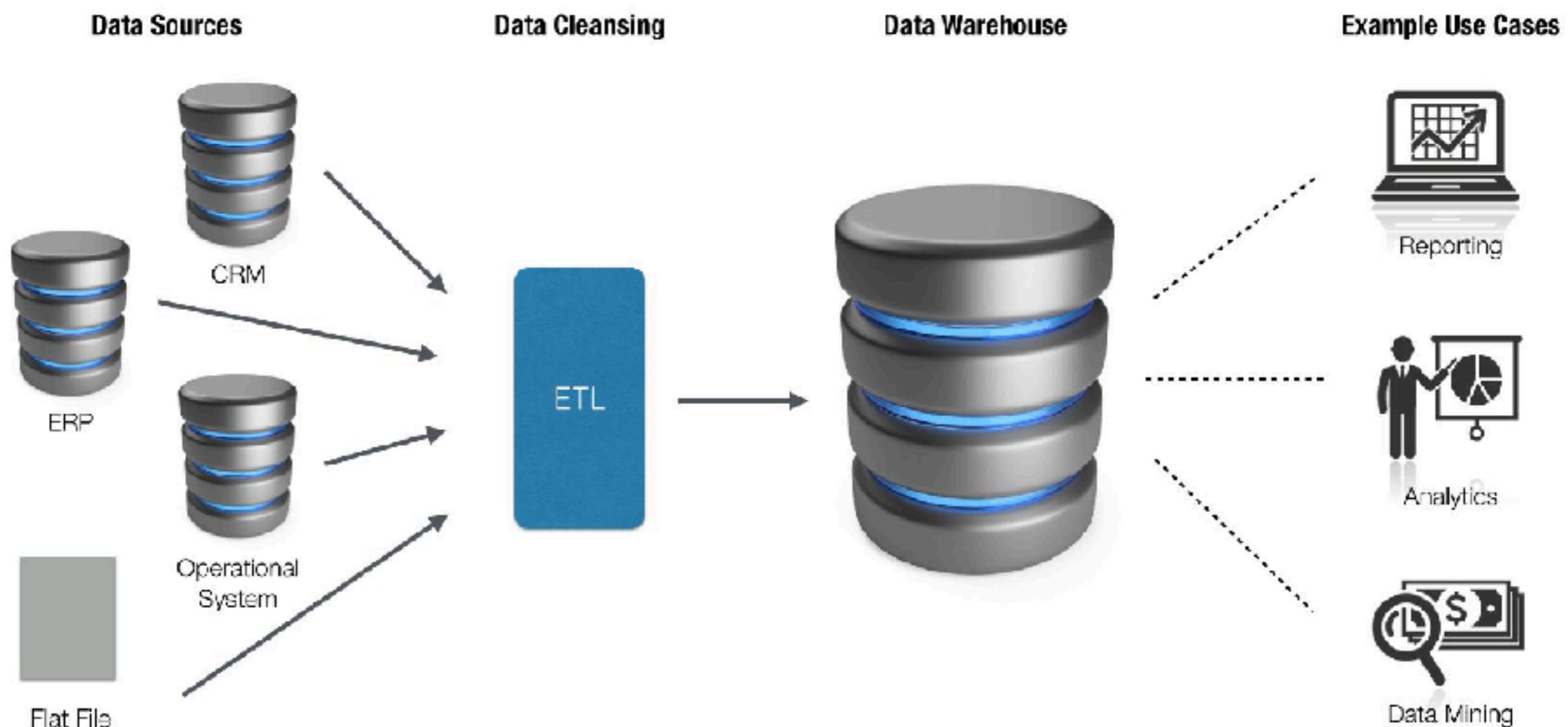


Source: Data Science and Critical Thinking, A. Croll

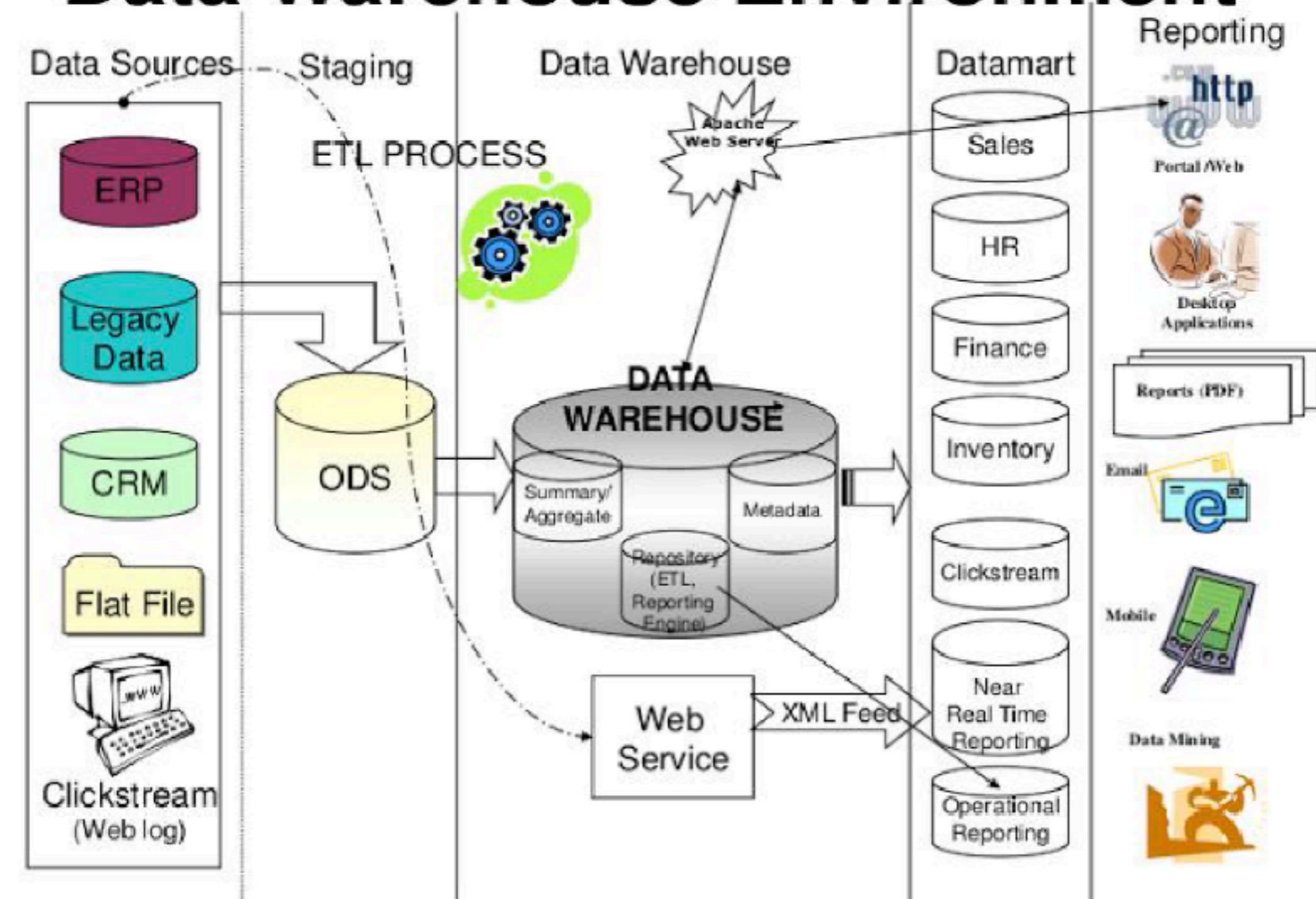


# Data Warehouse

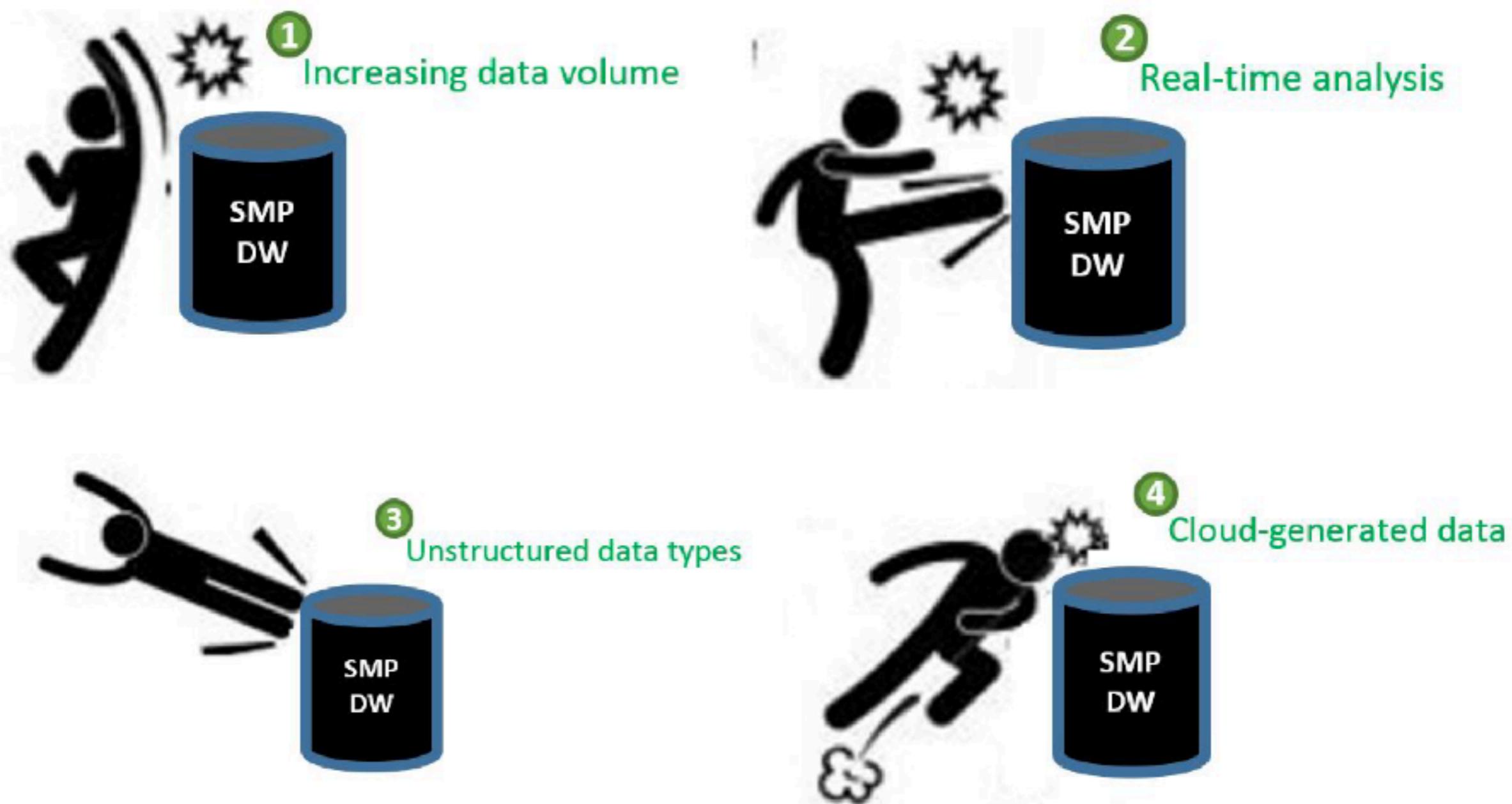
Traditional approach to integrating data for consistency and quality



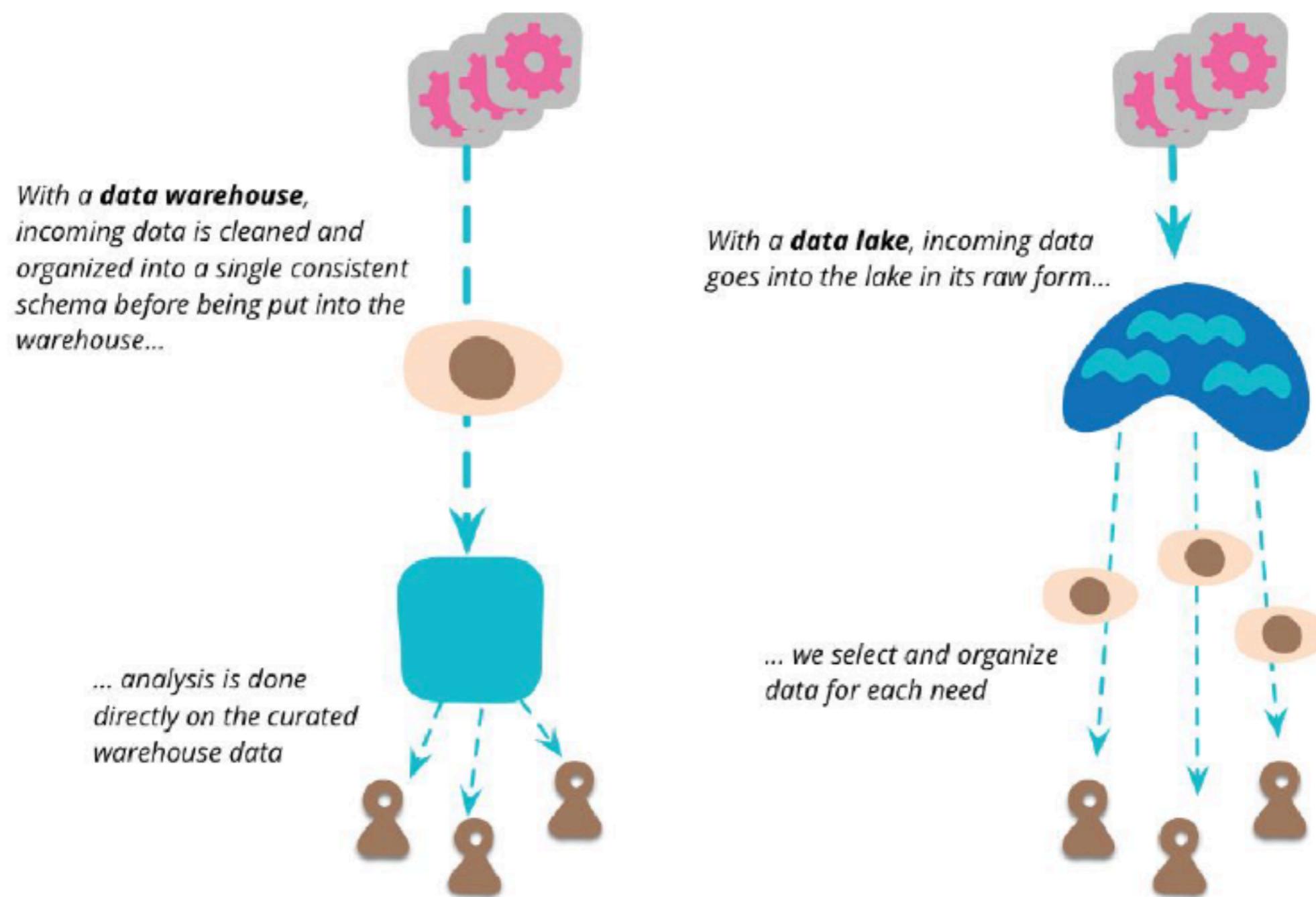
# Data Warehouse Environment

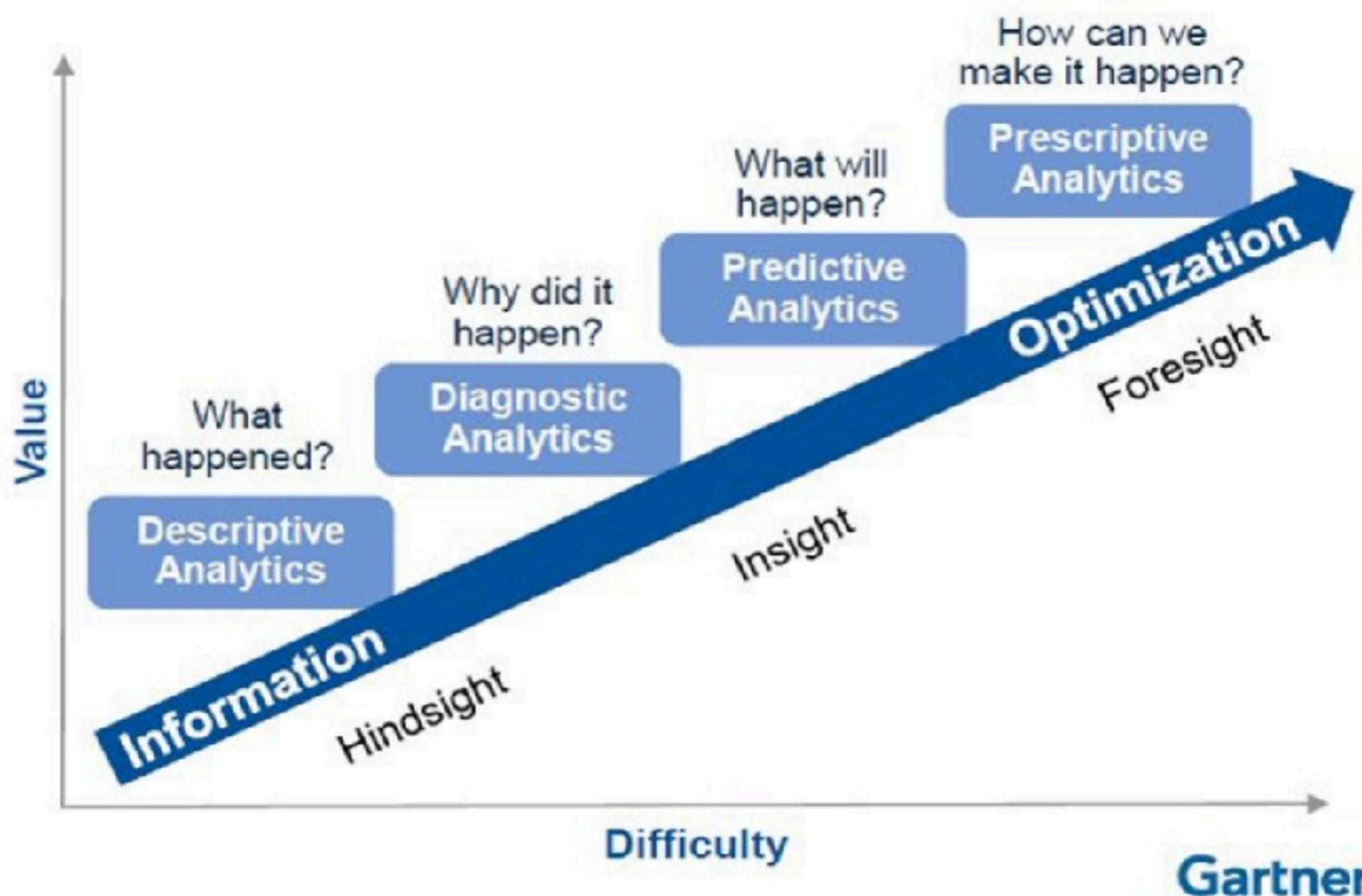


# Data Warehouse



# Differences between Data Lake and Data Warehouse





## Traditional Analytics (BI)

## vs

## Big Data Analytics

### Focus on

- Descriptive analytics
- Diagnosis analytics

### Data Sets

- Limited data sets
- Cleansed data
- Simple models

- **Predictive analytics**
- **Data Science**

### Supports

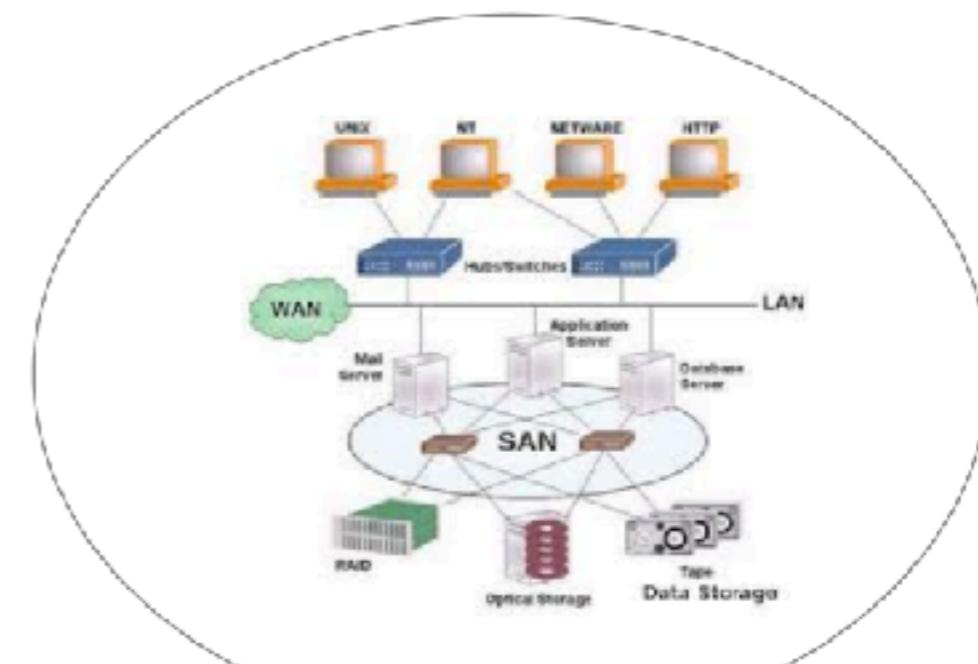
**Causation:** what happened,  
and why?

- Large scale data sets
- More types of data
- Raw data
- Complex data models

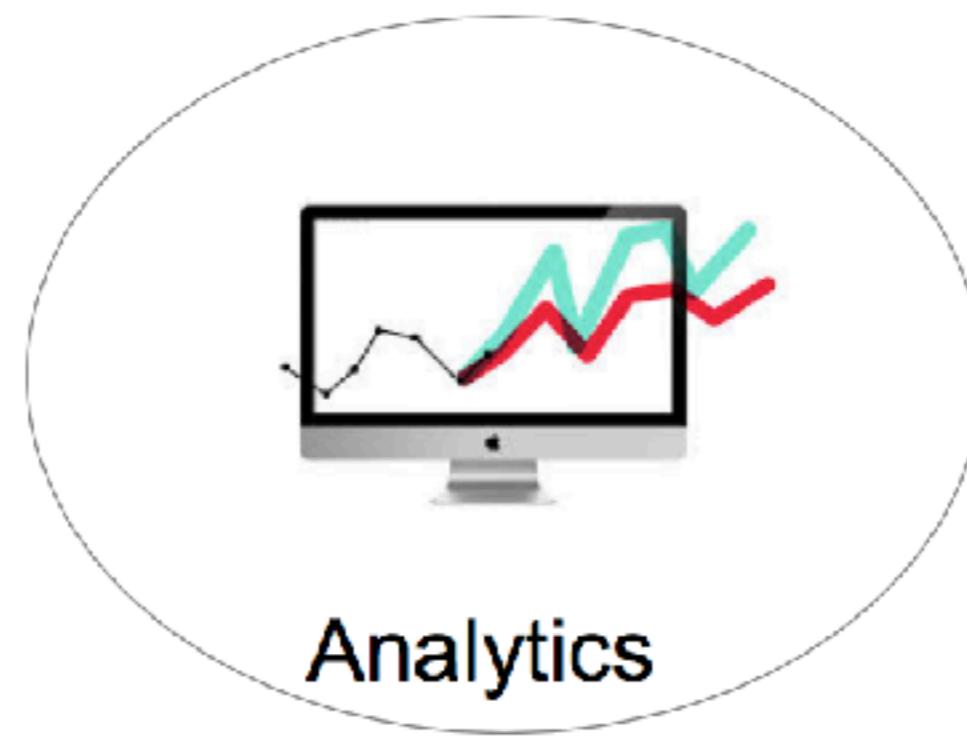
**Correlation:** new insight  
More accurate answers



Data Sources

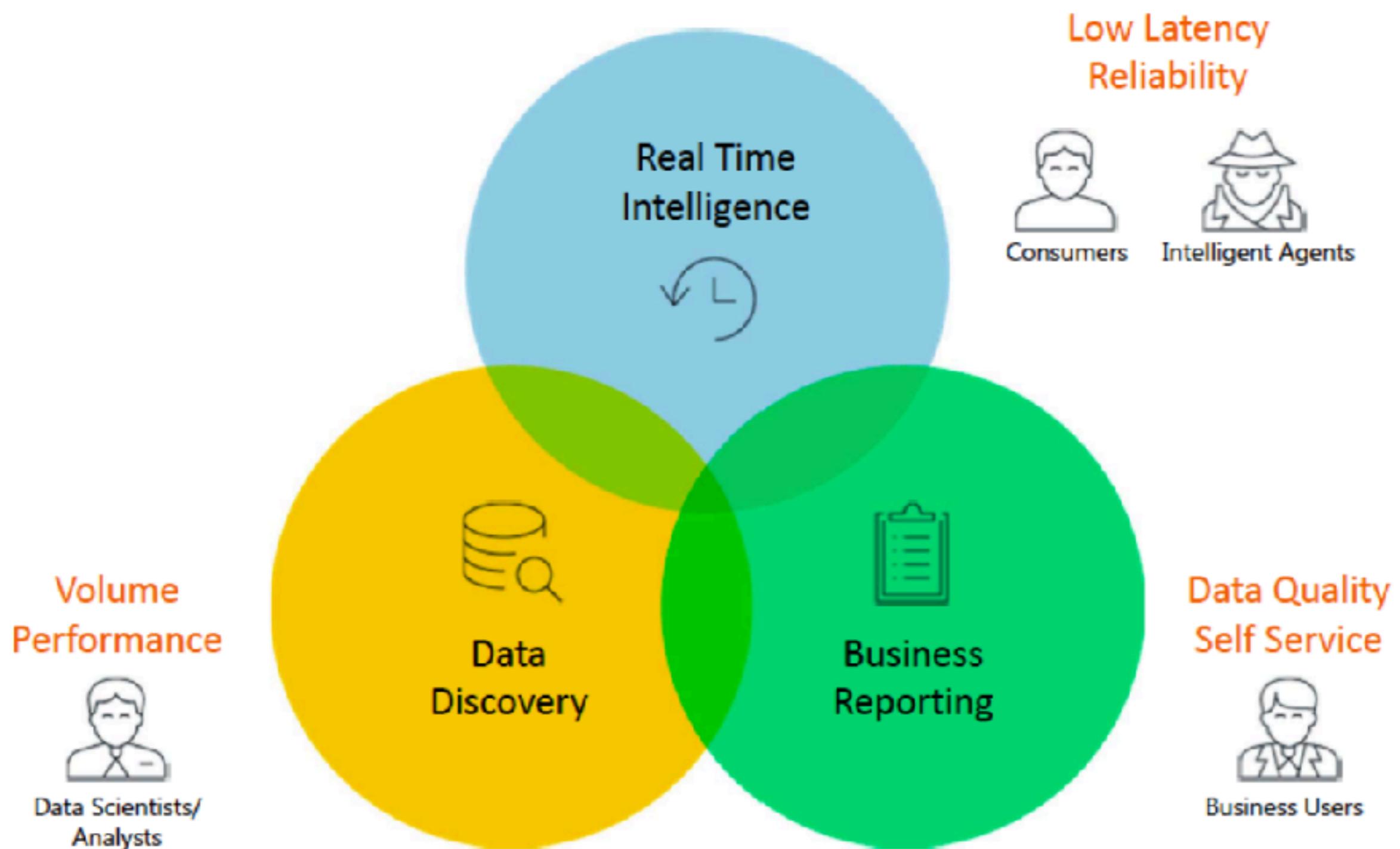


Technology

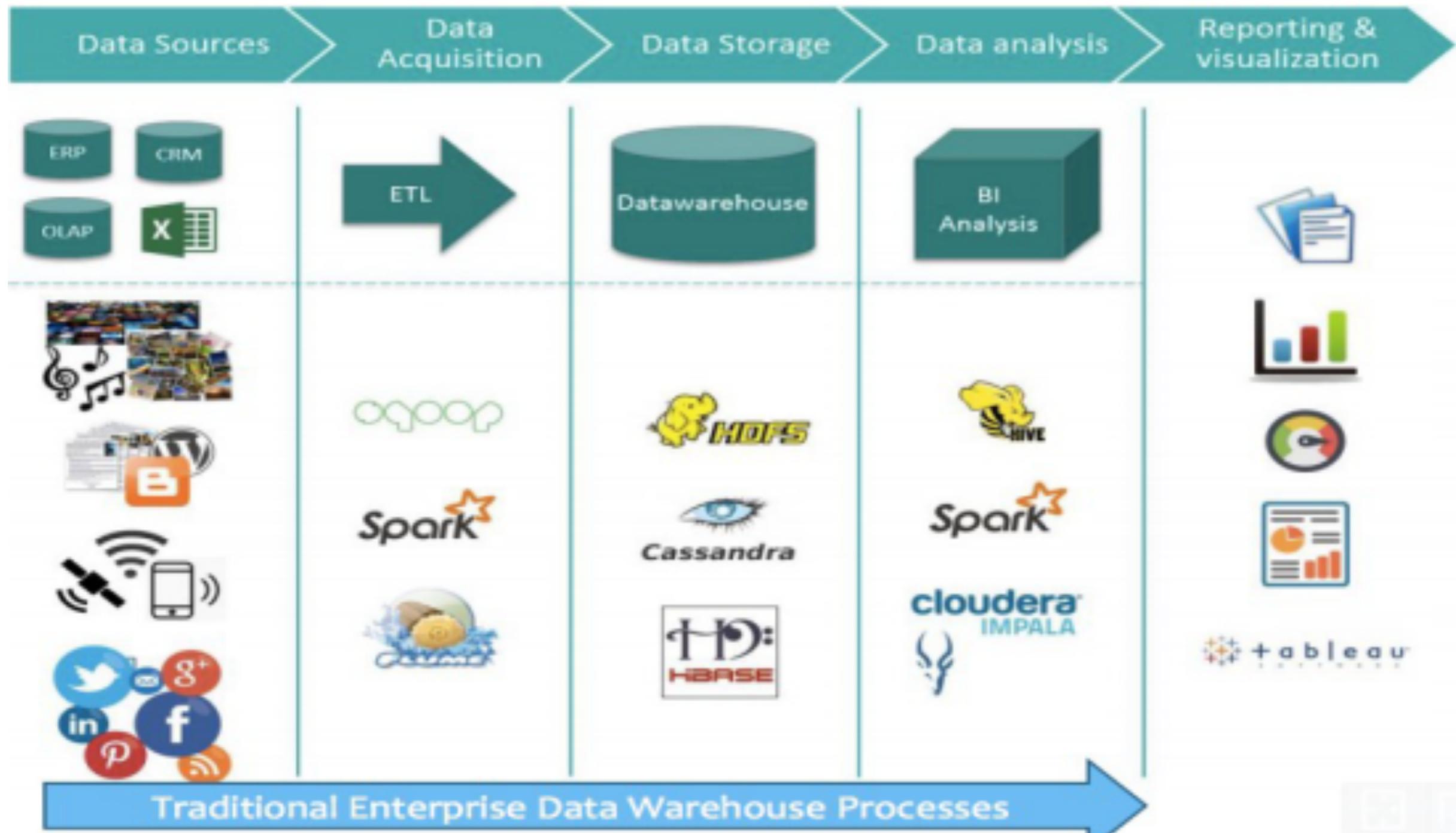


Analytics

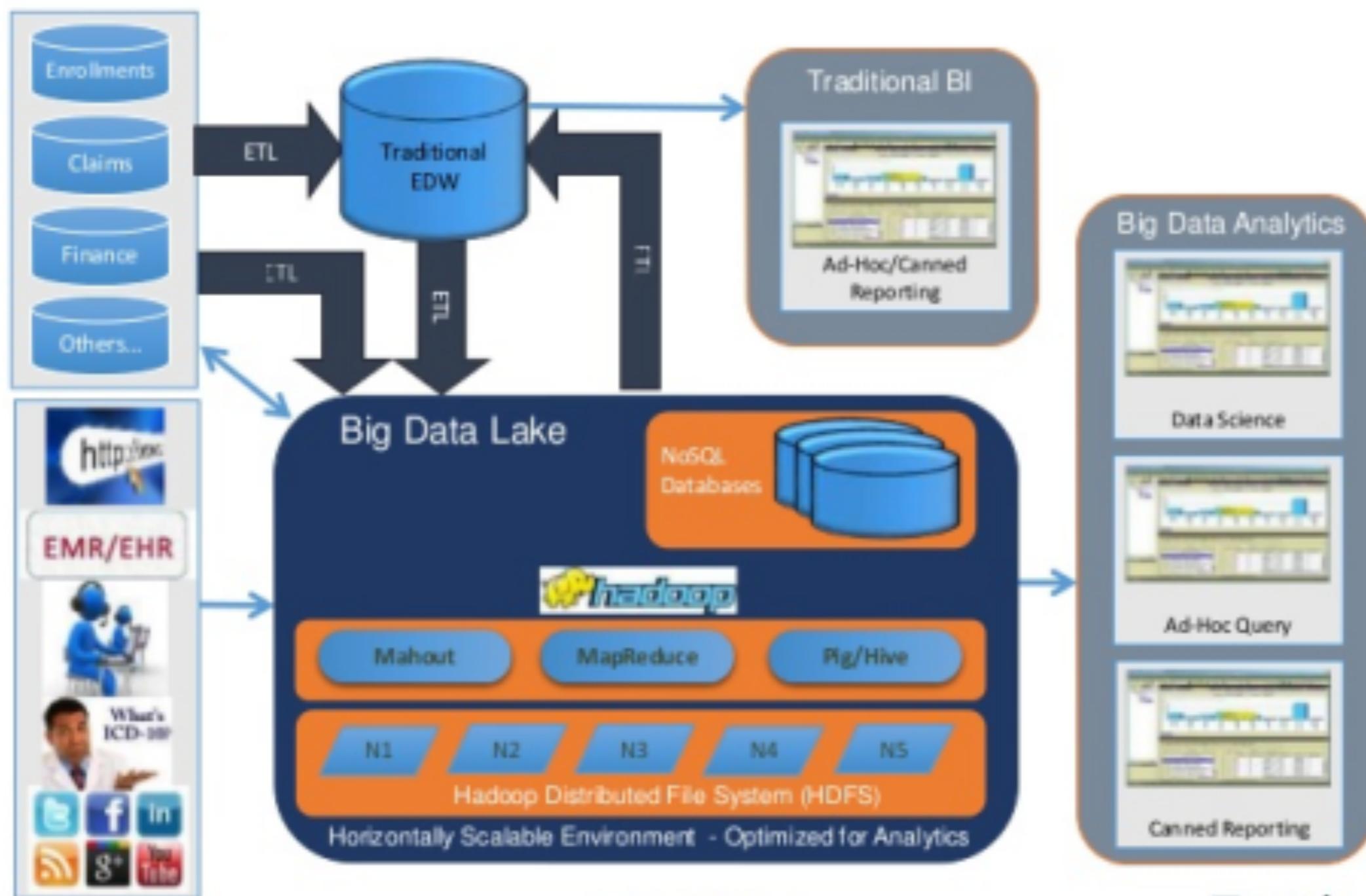
# Big Data Analytics



# How Data Lake Works?



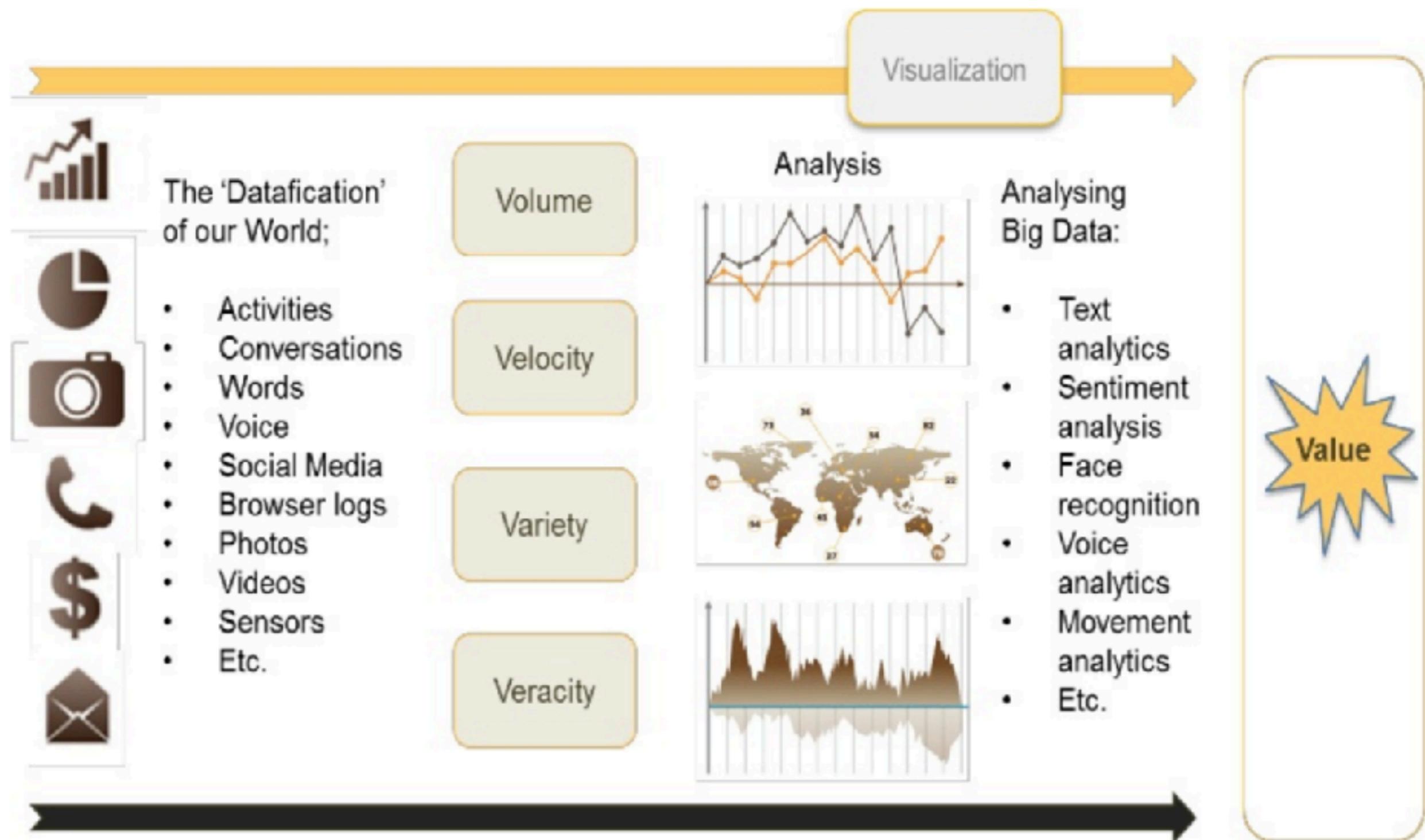
## Today's business environment requires Big Data



#edwdc15

@joe\_Caserta

 caserta



**Data Lake** isn't just a technology  
It is an architecture

# Data Lake: Key Benefits



- Scale as much as you can**
- Plug in disparate data sources**
- Acquire high-velocity data: Store in native format**
- Don't worry about schema**
- Unleash your favorite SQL**
- Advanced algorithms**
- Administrative resources**

# Data Lake v.s. Data WareHouse

Complementary to EDW (not replacement)	Data lake can be source for EDW
Schema on read (no predefined schemas)	Schema on write (predefined schemas)
Structured/semi-structured/Unstructured data	Structured data only
Fast ingestion of new data/content	Time consuming to introduce new content
Data Science + Prediction/Advanced Analytics + BI use cases	BI use cases only (no prediction/advanced analytics)
Data at low level of detail/granularity	Data at summary/aggregated level of detail
Loosely defined SLAs	Tight SLAs (production schedules)
Flexibility in tools (open source/tools for advanced analytics)	Limited flexibility in tools (SQL only)

**More Data Sources**

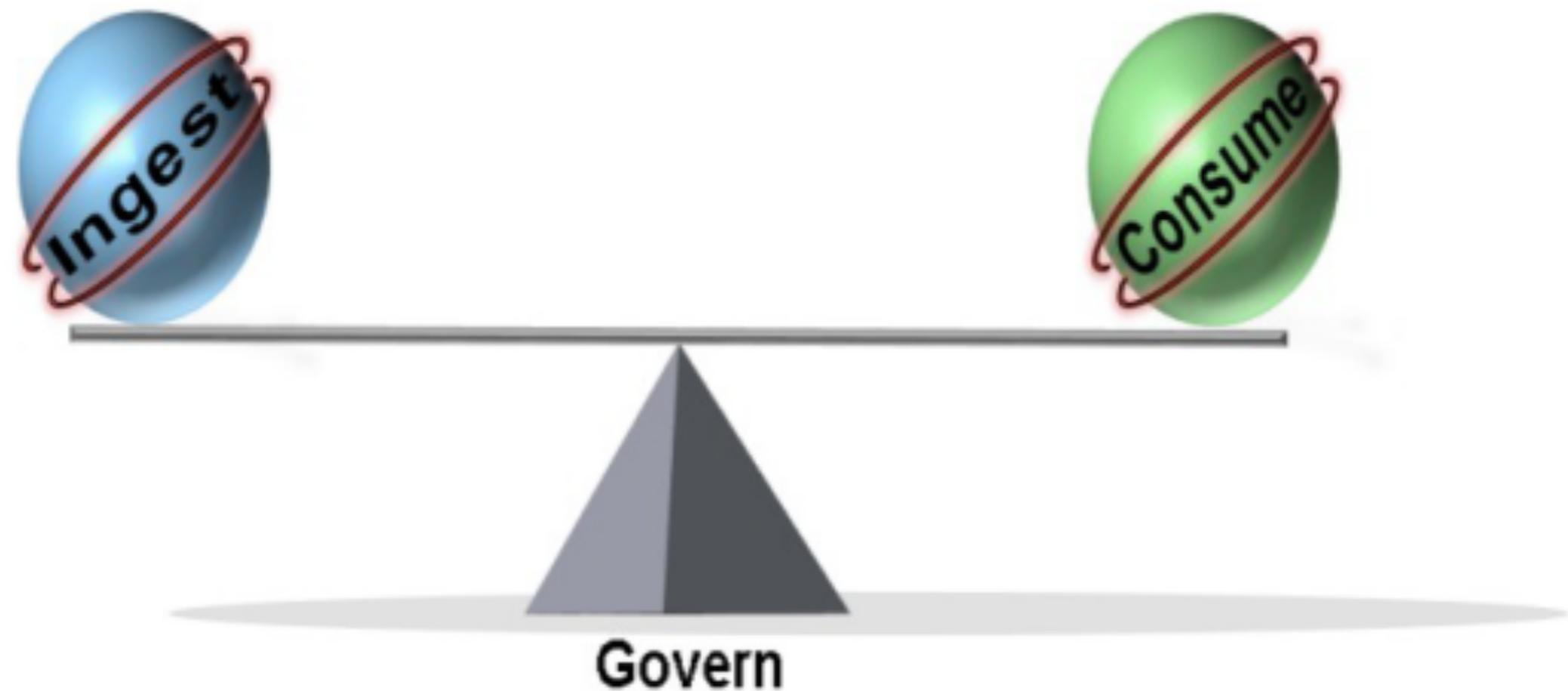
**More Applications**

**More Business Units**

**More Users**

**Without proper governance mechanisms  
Data lakes risk turning data swamps**

# Data Lake Governance



Source: What is “Just-Enough” Governance for the Data Lake?

## Fundamental Capabilities

- **The definition of the incoming data from a Business use perspective**
- **Documentation of the context, lineage, and frequency of the incoming data;**
- **Security level classification of the incoming data;**
- **Documentation of creation, usage, privacy, regulatory, and encryption business rules which apply to the incoming data.**

# What can it do for my Data Lake

- **Where did my data come from ? How is it being transformed ?**
- **Track usage, resolve anomalies, visualize, optimize and clarify data lineage**
- **Search and access data**
- **Assess data quality and fitness for purpose**
- **Govern who can/cannot access the data**
- **Data life cycle management, archiving and retention policies**
- **Auditing, compliance**

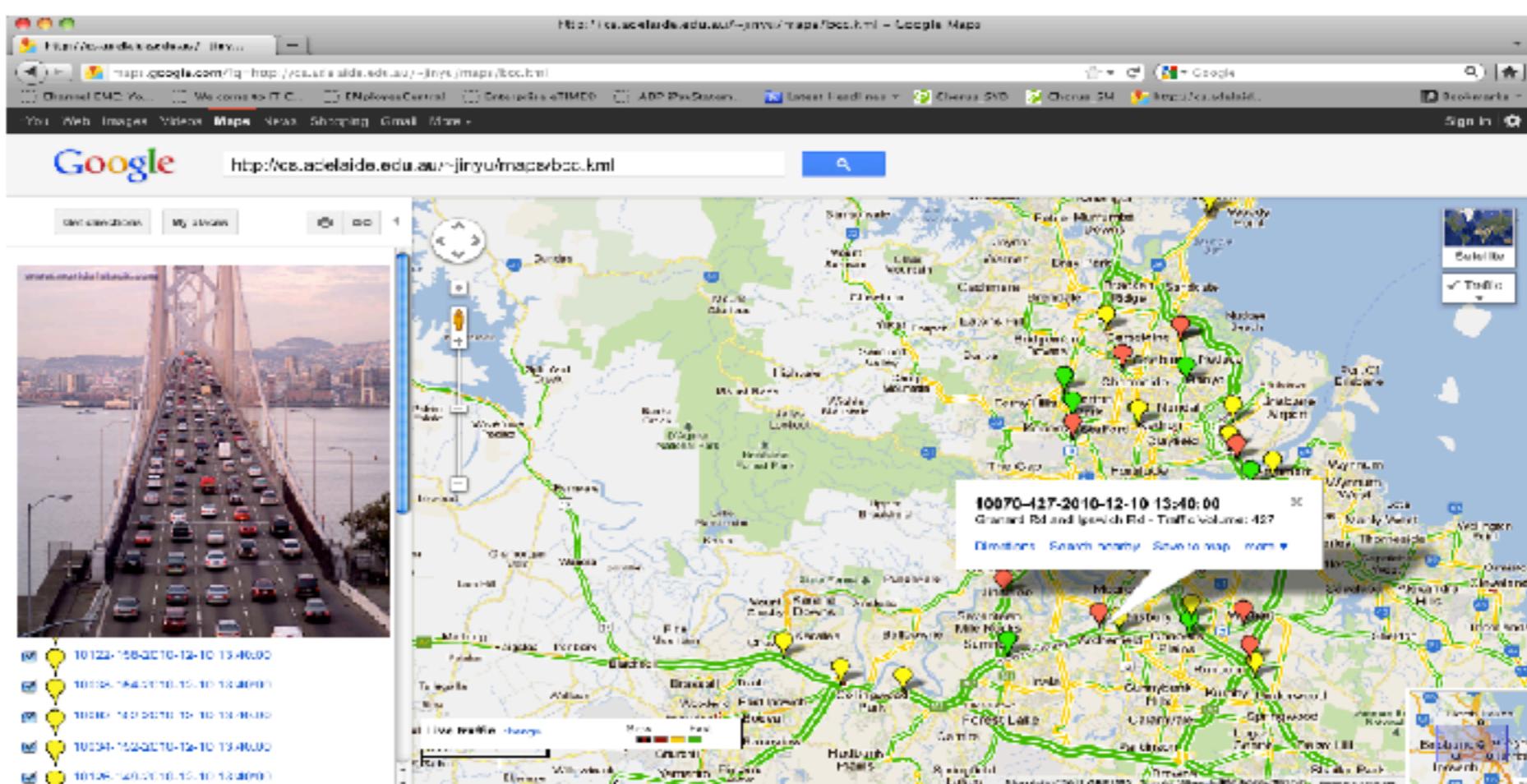
# Summary

- **Big Data: Data lake instead of data warehouse?**
- **Data Lake is not only a technology, it is an architecture**
- **Data Lake components: Data acquisition (intake), Data management, Data Storage, Data consumption (Discovery)**
- **Data Lake governance is very important**

## Customer Example: Analytics

### Municipal Traffic Analysis to Simulate Traffic Velocity Patterns and Reduce Delays

- Correlate multiple types of data (GPS, weather, sensor, video, social media)
- Simulation techniques to model traffic transition points
- Signal retiming to minimize stops and delays
- Peak delays reduced by 16% and stops reduced by 22%



Pivotal

# Customer Case Studies on Big Data Lake

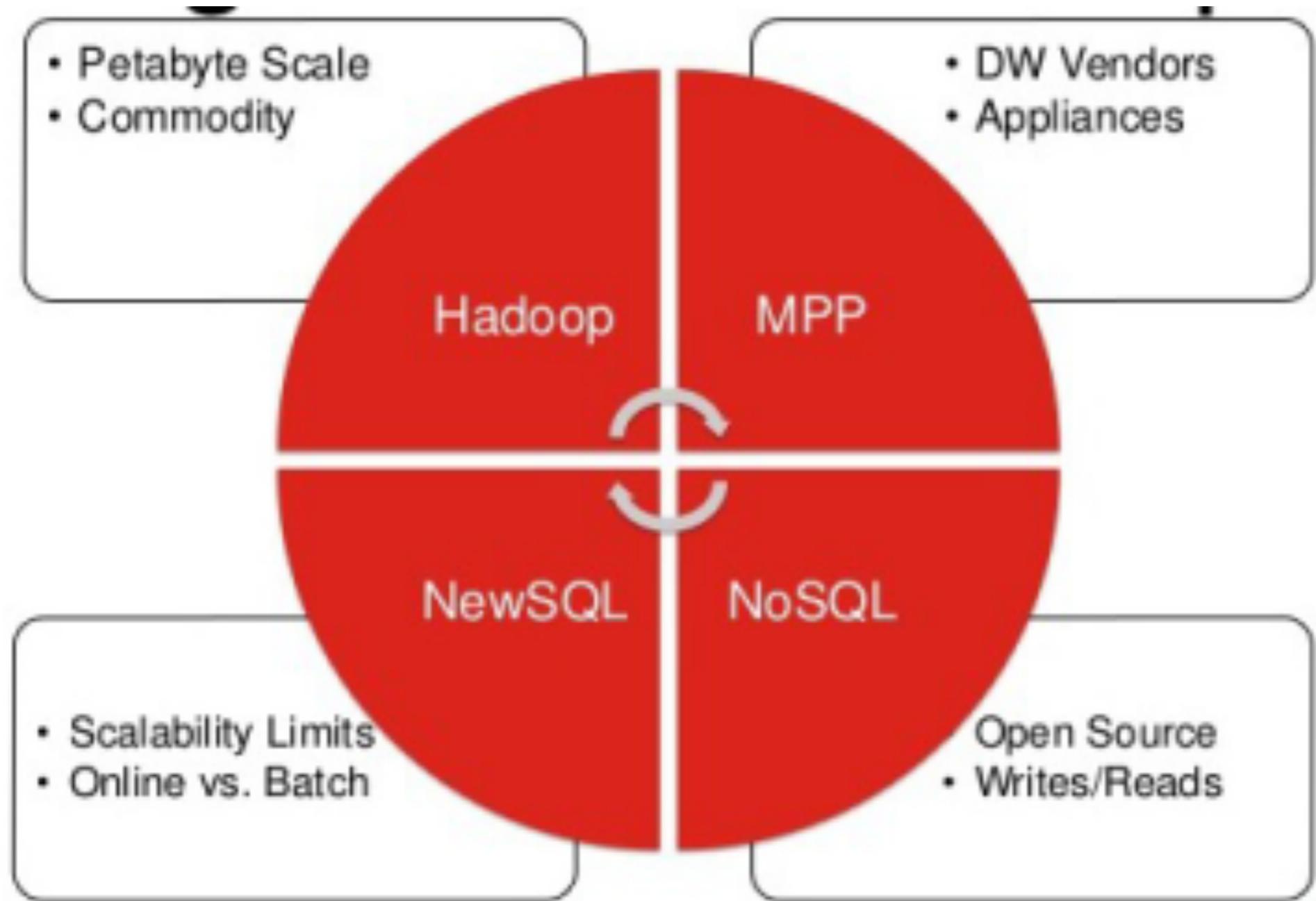
## China Railways scales online sales for the largest rail way in the world with Pivotal Gemfire

- Reliable, High Performance and Continuous uptime with thousands of transactions per second
- On demand scaling for growth
- Cost effective operations



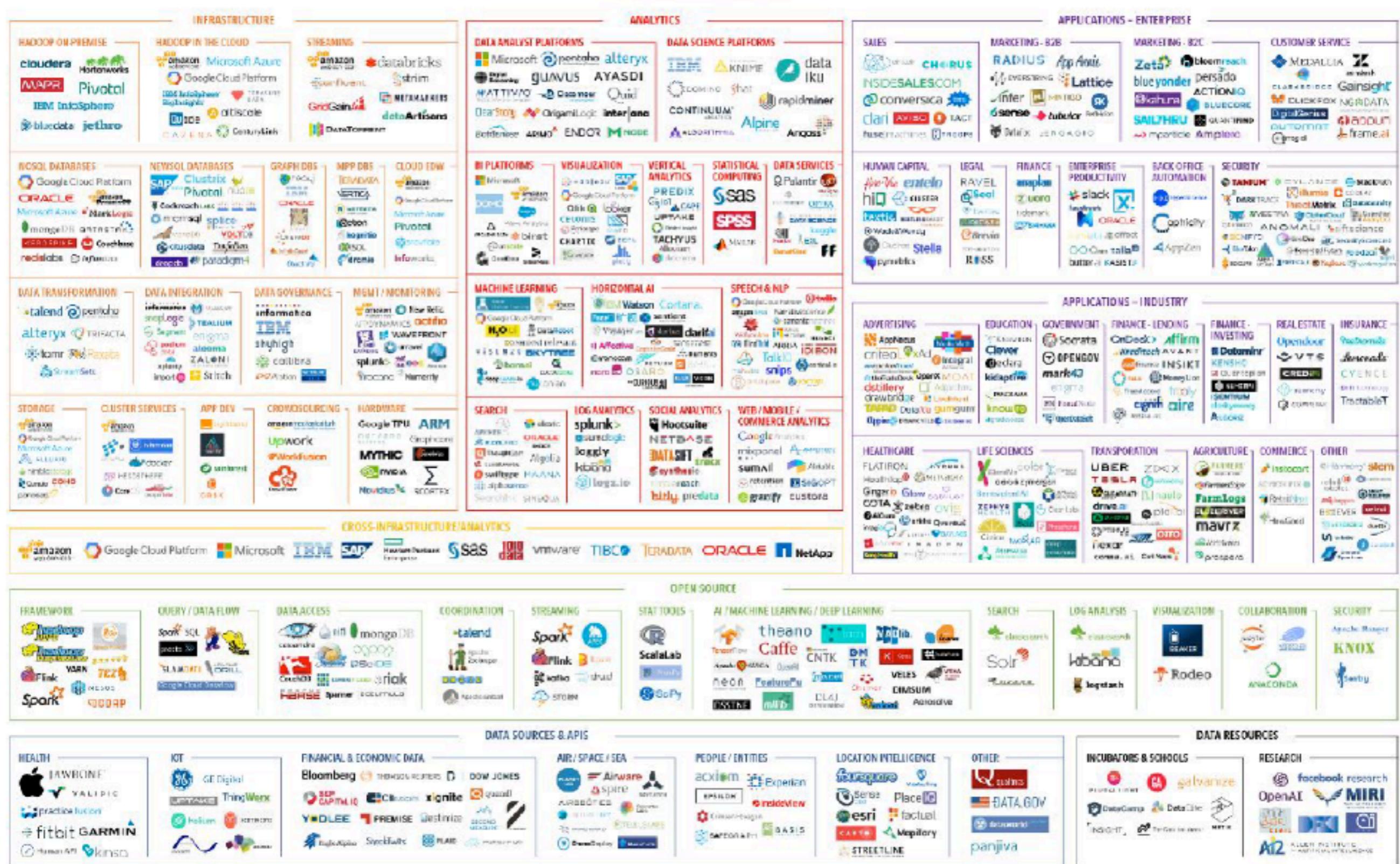
# Big Data Technology

# Big Data Landscape

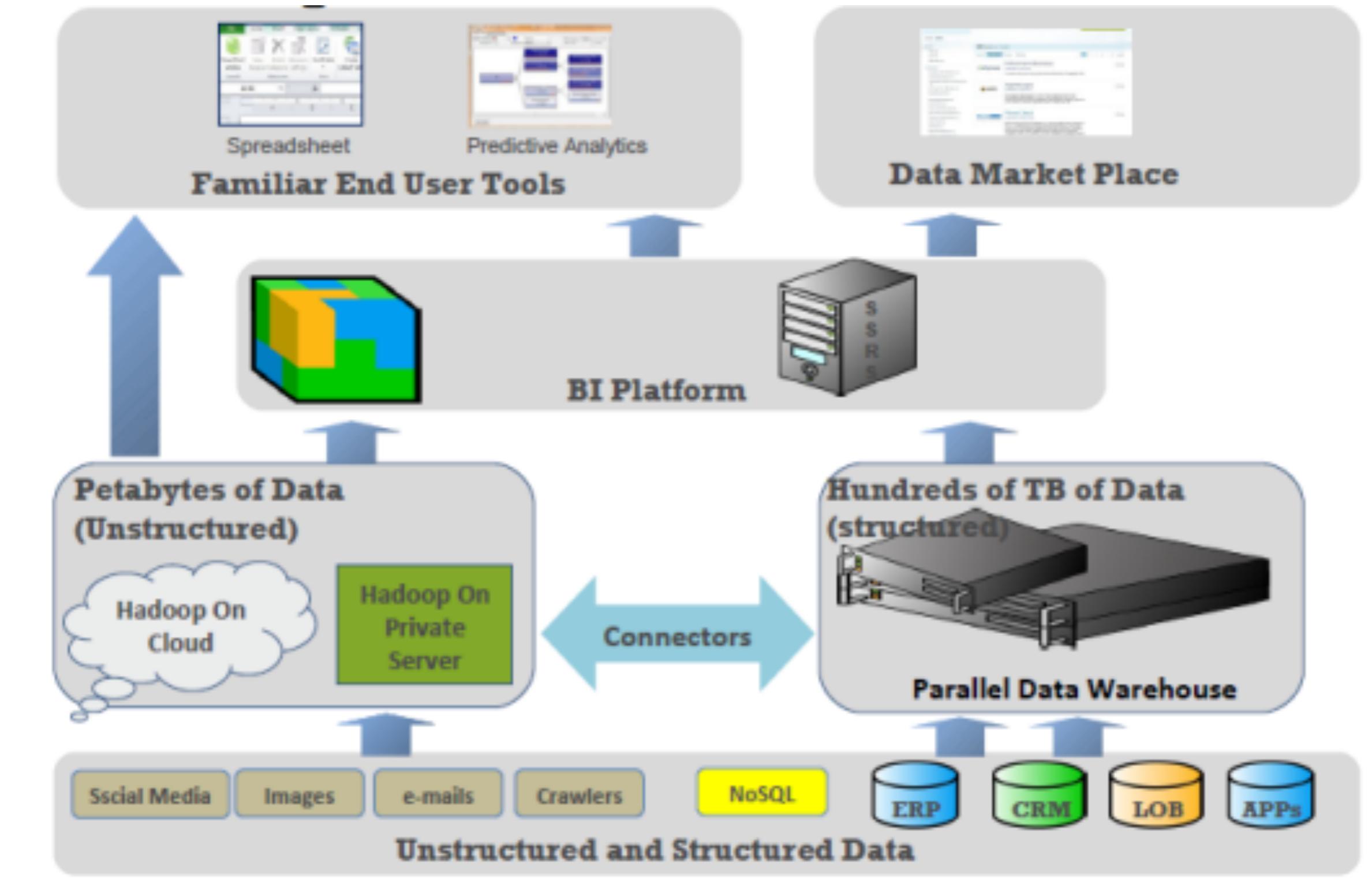


Source: Big Data in the Enterprise. When to Use What?

# Big Data Landscape 2017



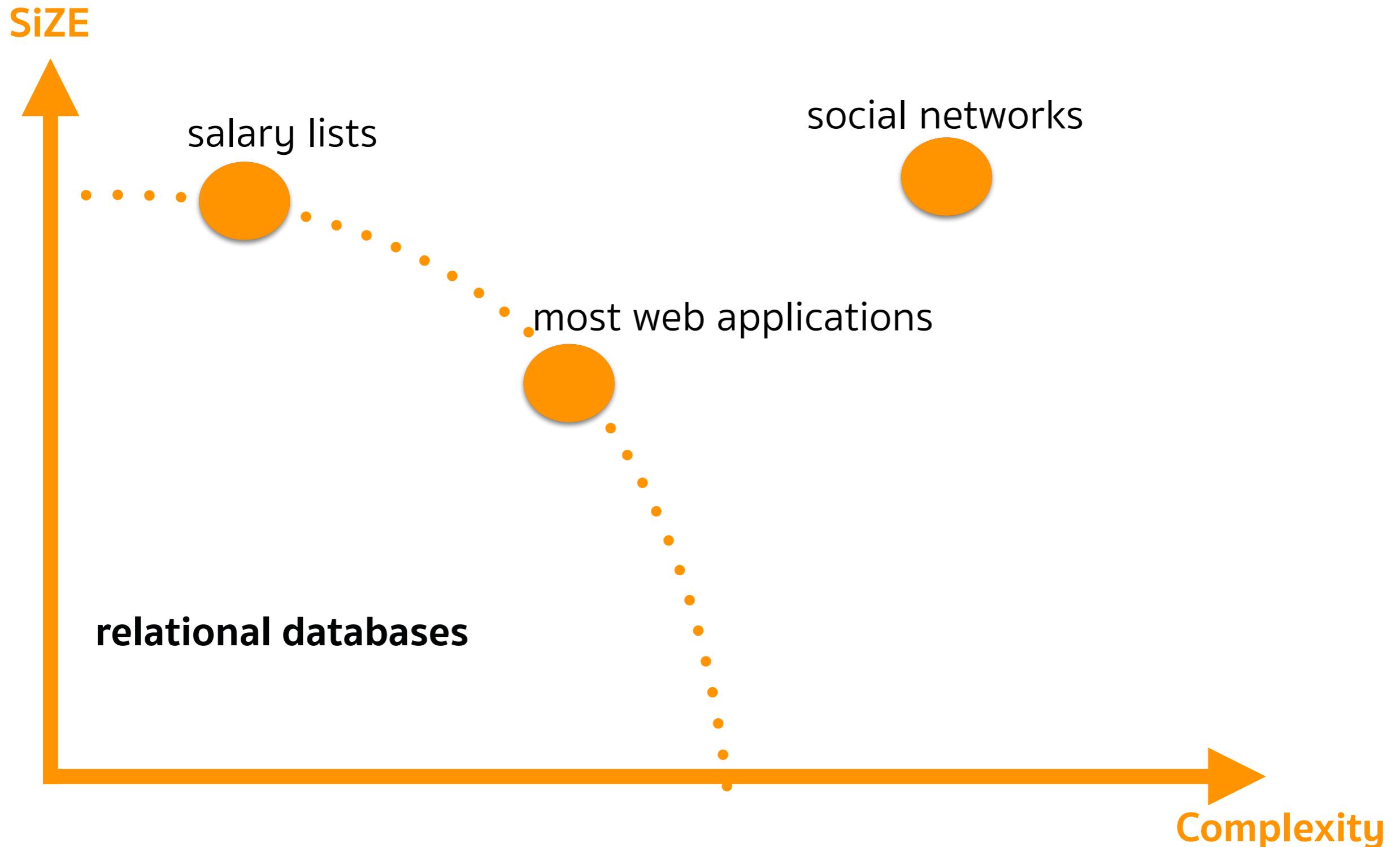
# Big Data Future Architecture



# What is NoSQL ?

A NoSQL (Not only SQL) database provides a mechanism for storage and retrieval of data that is modeled in means other than the tabular relations used in RDBMS.

Motivations for this approach include simplicity of design, horizontal scaling, and finer control over availability.



# NoSQL PROS AND CONS

## PROS

**MASSIVE SCALABILITY**

**HIGH AVAILABILITY**

**LOWER COST**

**SCHEMA FLEXIBILITY**

**STRUCTURED AND SEMI STRUCTURED DATA**

## CONS

**LIMITED QUERY CAPABILITIES**

**NOT STANDARDISED (PORTABILITY MAY BE AN ISSUE)**

**STILL A DEVELOPING TECHNOLOGY**

# Types of NoSQL

**Column-oriented**

**Key Value Store**

**Document Store**

**Graph**

# Column-oriented databases

Row Oriented  
(RDBMS Model)

<b>id</b>	<b>Name</b>	<b>Age</b>	<b>Interests</b>
1	Ricky		Soccer, Movies, Baseball
2	Ankur	20	
3	Sam	25	Music

Multi-valued

null

Column Oriented  
(Multi-value sorted map)

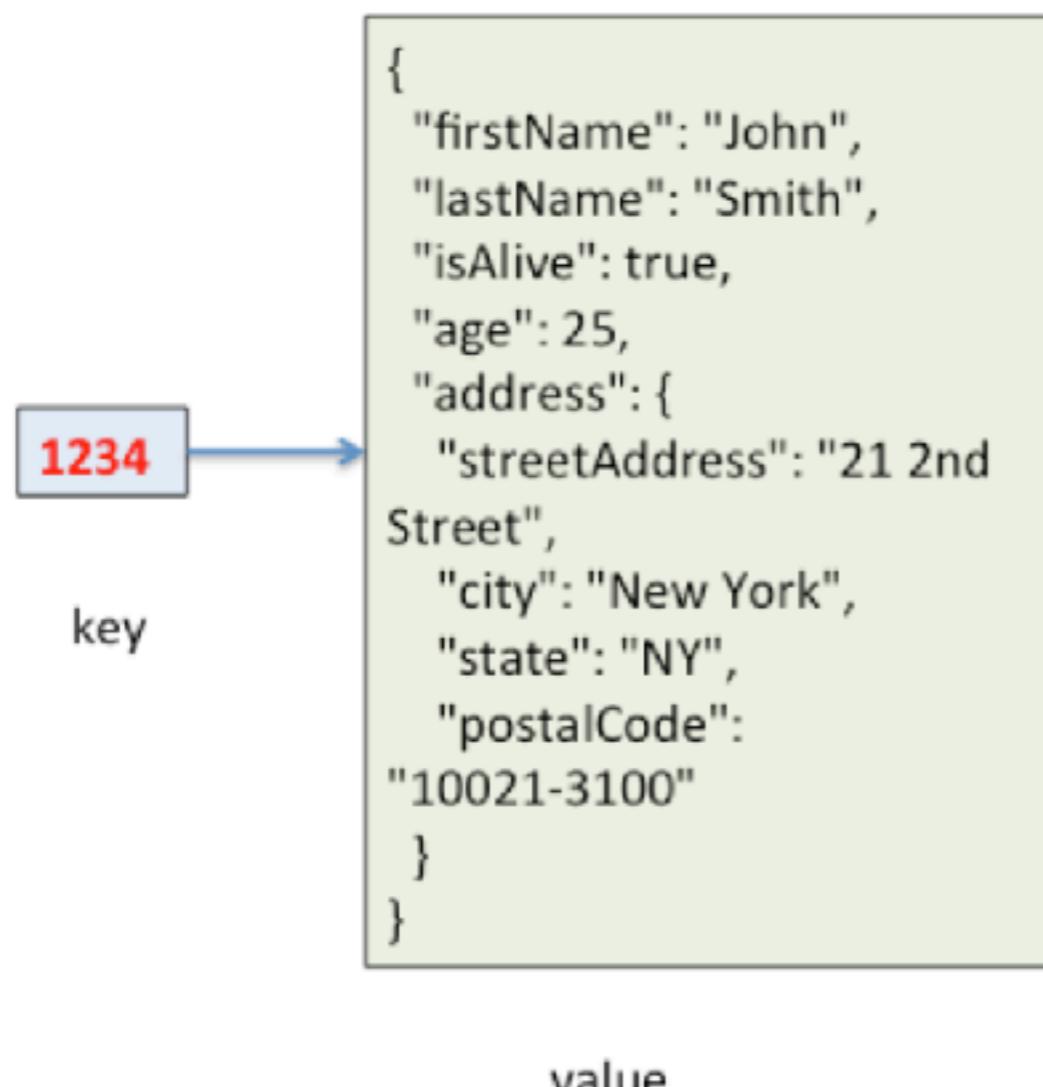
<b>id</b>	<b>Name</b>	<b>id</b>	<b>Age</b>	<b>id</b>	<b>Interests</b>
1	Ricky	2	20	1	Soccer
2	Ankur	3	25	1	Movies
3	Sam			1	Baseball
				3	Music



# Key-value store database

The storage of a value against a key

A key-value store requires the key to be specified and the key must be known to retrieve the value



Key	Value
Mahesh	{"Mathematics, Science, History, Geography"}
Uma	{"English, Hindi, French, German"}
Paul	{"Computers, Programming"}
Abraham	{"Geology, Metallurgy, Material Science"}



redis



# Document-oriented database

**Designed for storing, retrieving, and managing document-oriented information, also known as semi-structured data.**

**Most of the databases available under this category use**

**XML, JSON, BSON, or YAML**

```
{  
    "EmployeeID": "SM1",  
    "FirstName" : "Anuj",  
    "LastName"  : "Sharma",  
    "Age"        : 45,  
    "Salary"     : 10000000  
}  
  
{  
    "EmployeeID": "MM2",  
    "FirstName" : "Anand",  
    "Age"        : 34,  
    "Salary"     : 5000000,  
    "Address"   : {  
        "Line1"  : "123, 4th Street",  
        "City"   : "Bangalore",  
        "State"  : "Karnataka"  
    },  
    "Projects" : [  
        "nosql-migration",  
        "top-secret-007"  
    ]  
}
```

# Document-oriented database



Comment Table  
Reader Table  
Article Table  
Author Table

Relational Database approach  
Document store approach

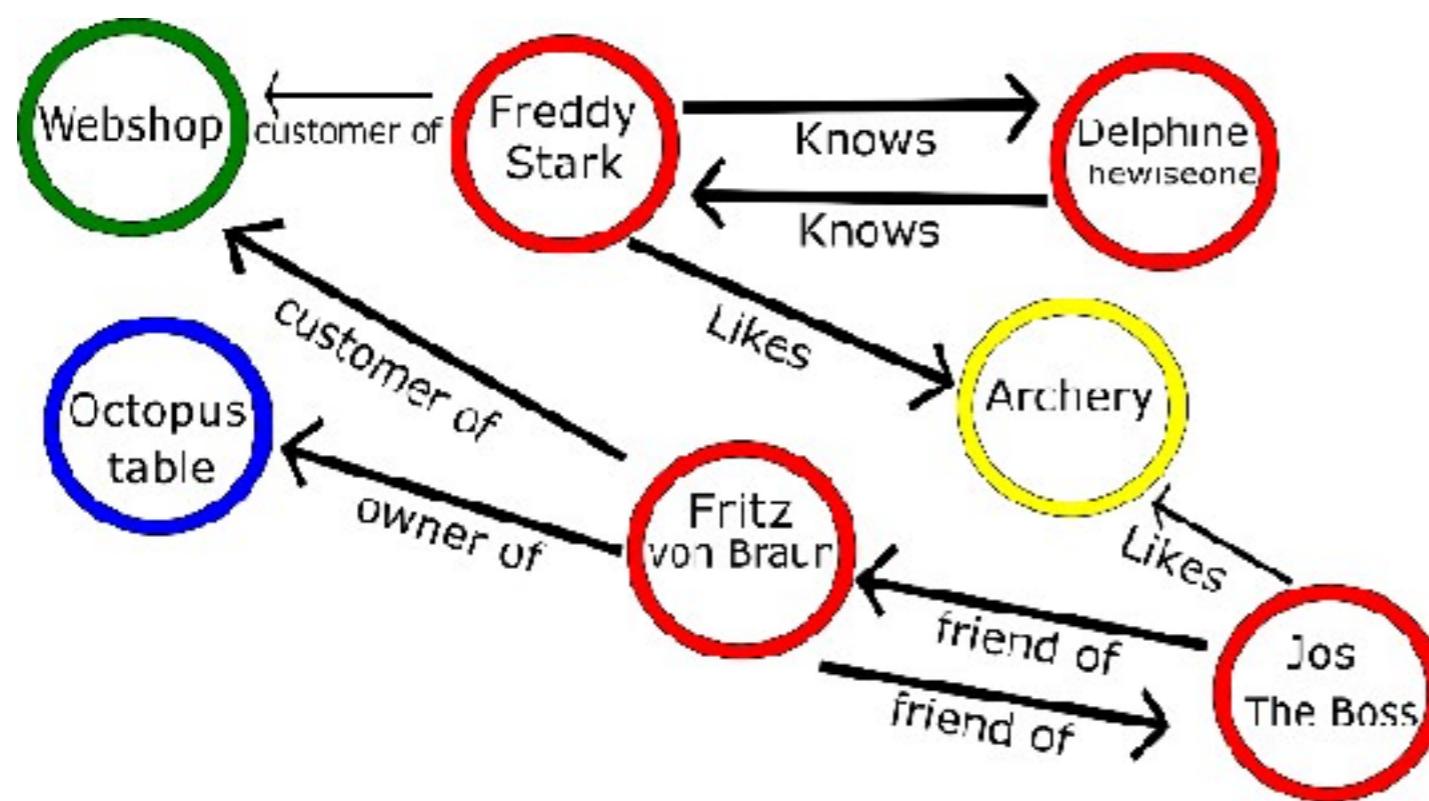
```
{  
  "articles": [  
    {  
      "title": "title of the article",  
      "articleID": 1,  
      "body": "body of the artricle",  
      "author": "Isaac Asimov",  
      "comments": [  
        {  
          "username": "Fritz",  
          "join date": "1/4/2014",  
          "commentid": 1,  
          "body": "this is a great article",  
          "replies": [  
            {  
              "username": "Freddy",  
              "join date": "11/12/2013",  
              "commentid": 2,  
              "body": "seriously? it's rubbish"  
            }  
          ]  
        },  
        {  
          "username": "Stark",  
          "join date": "19/06/2011",  
          "commentid": 3,  
          "body": "I don't agree with the conclusion"  
        }  
      ]  
    }  
  ]  
}
```

Whereas relational databases chop up data, Document stores save documents as a single entity



# Graph database

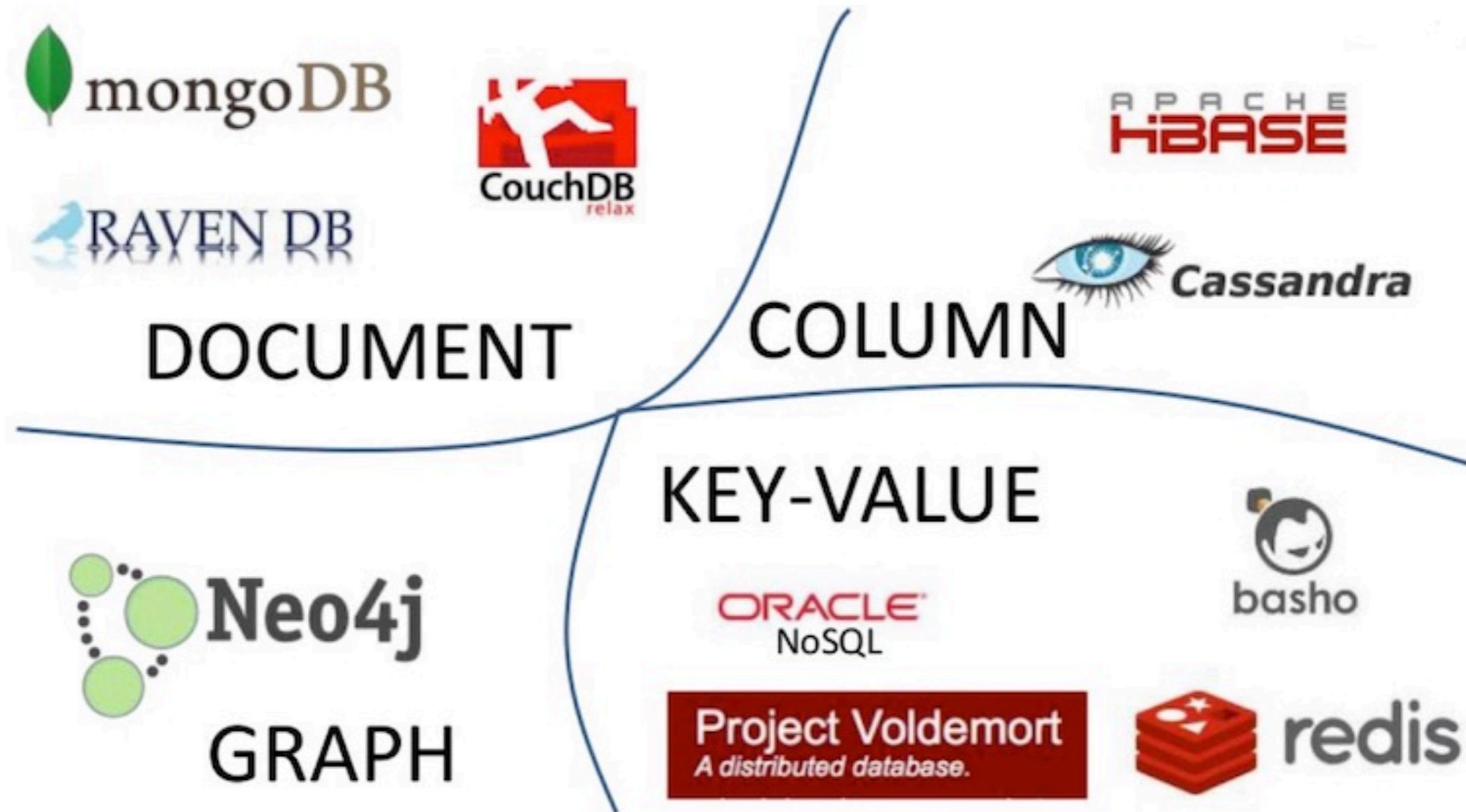
A database that uses graph structures for semantic queries with nodes, edges, and properties to represent and store data.

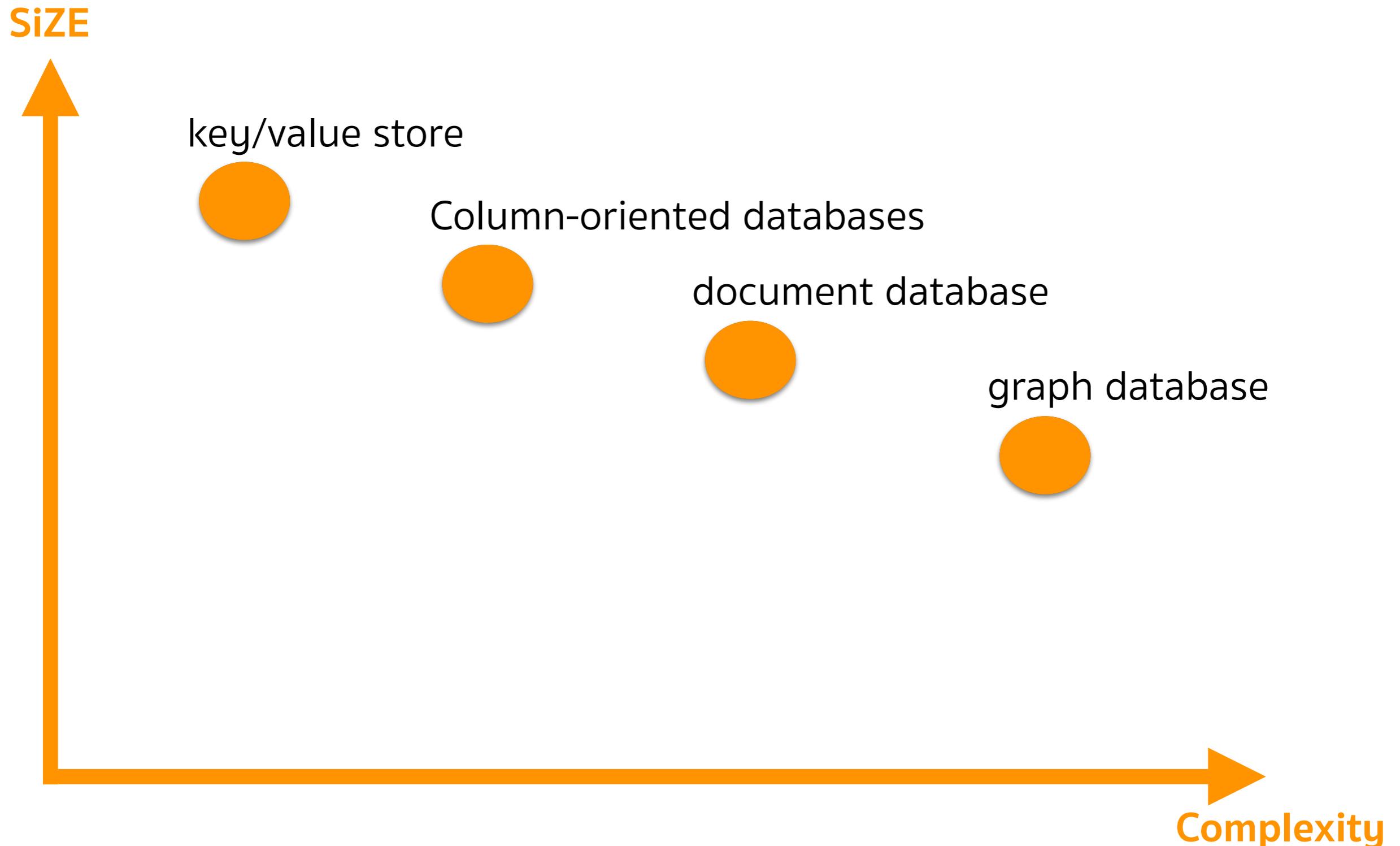


 Neo4j  
the graph database

 Franz Inc.  
**AllegroGraph**

 InfiniteGraph





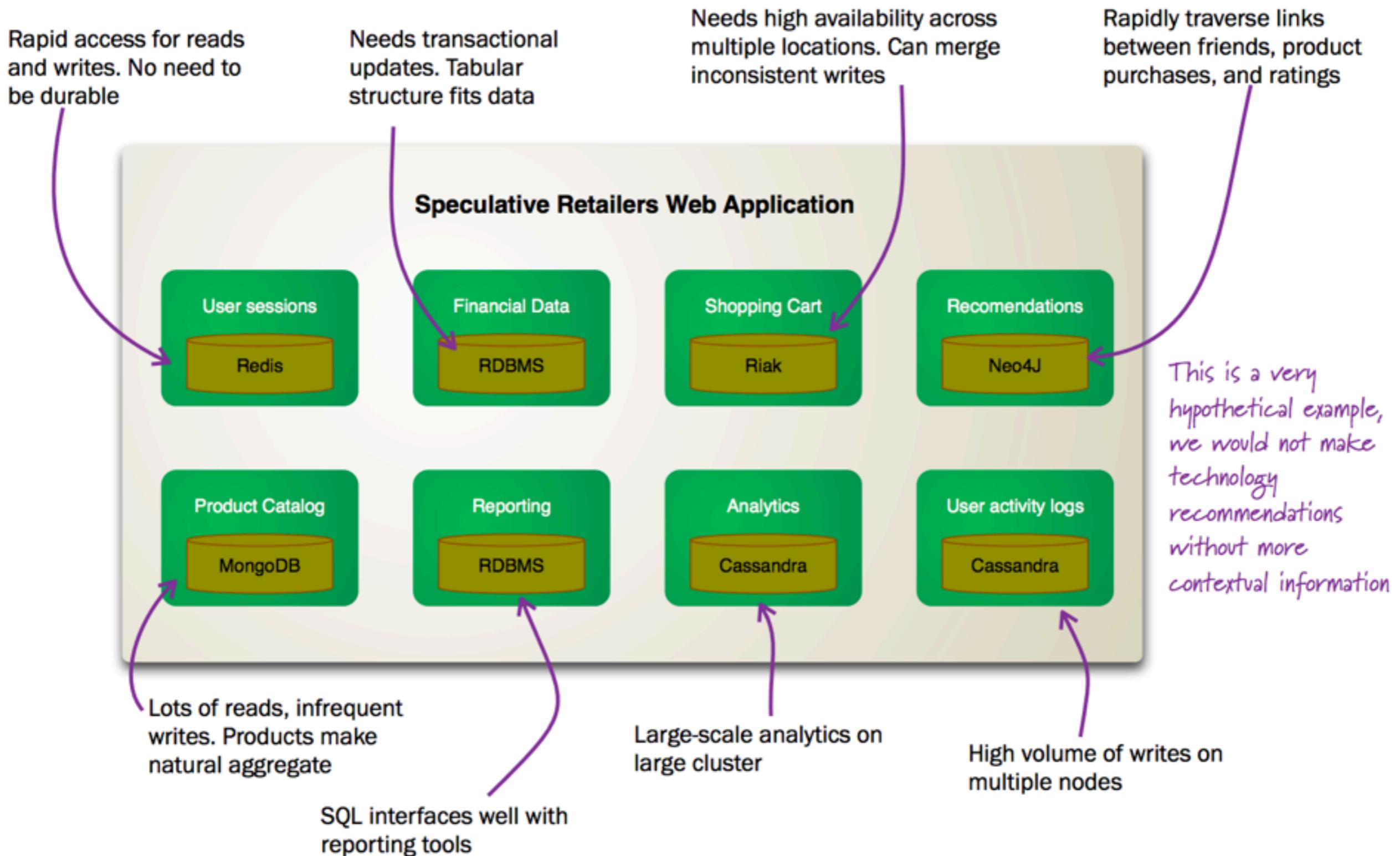
# SQ<sub>L</sub>

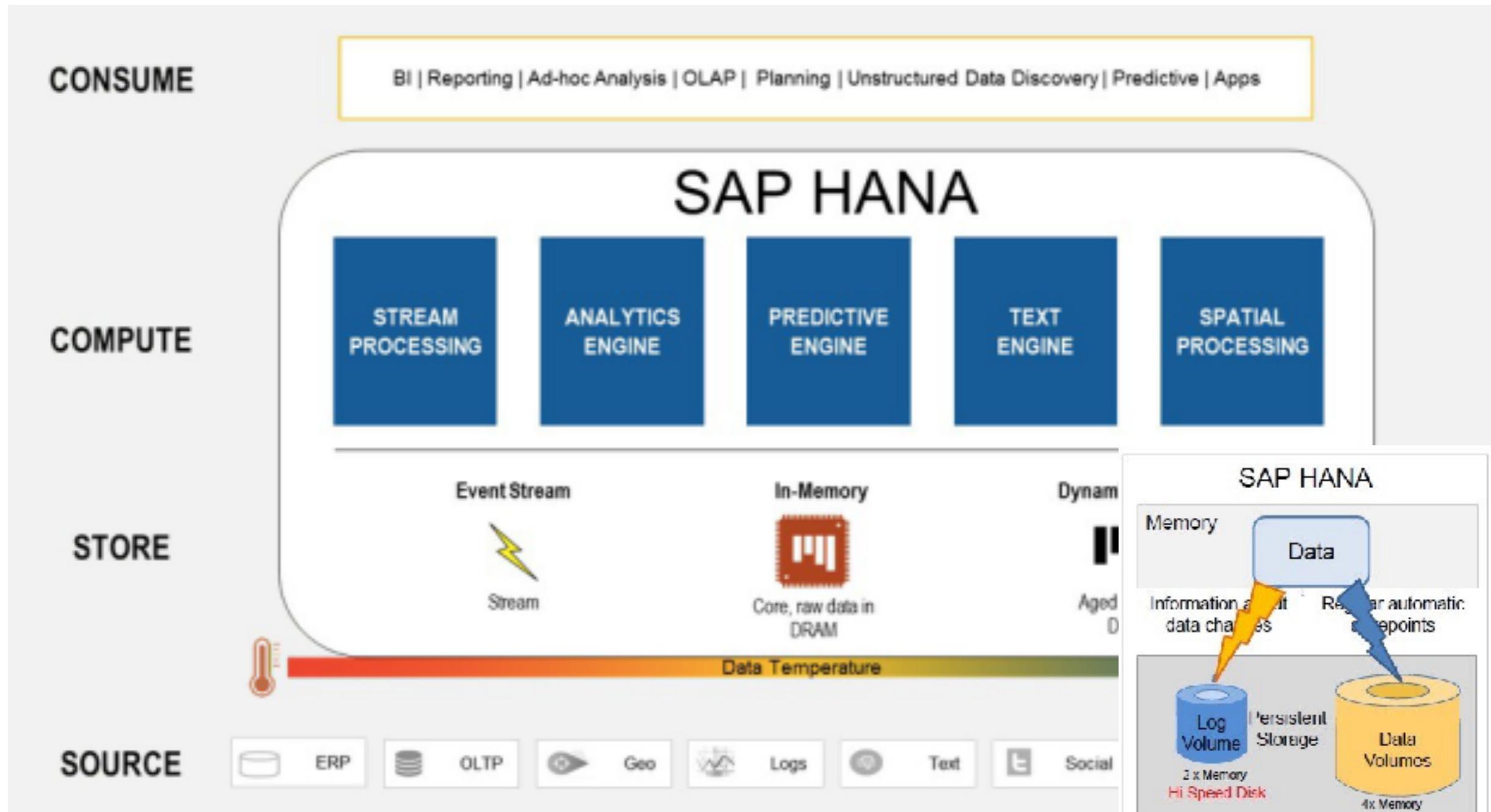
works great, can't scale for large data

# NoSQ<sub>L</sub>

works great, doesn't fit all situations

so use both, but think about when you want to use them!







**Oracle Exadata Database Machine**

## Extreme Performance for the Cloud

[Ellison announces next-generation systems](#) >

[Ease compliance: OFSAA and Oracle Exadata \(PDF\)](#) >

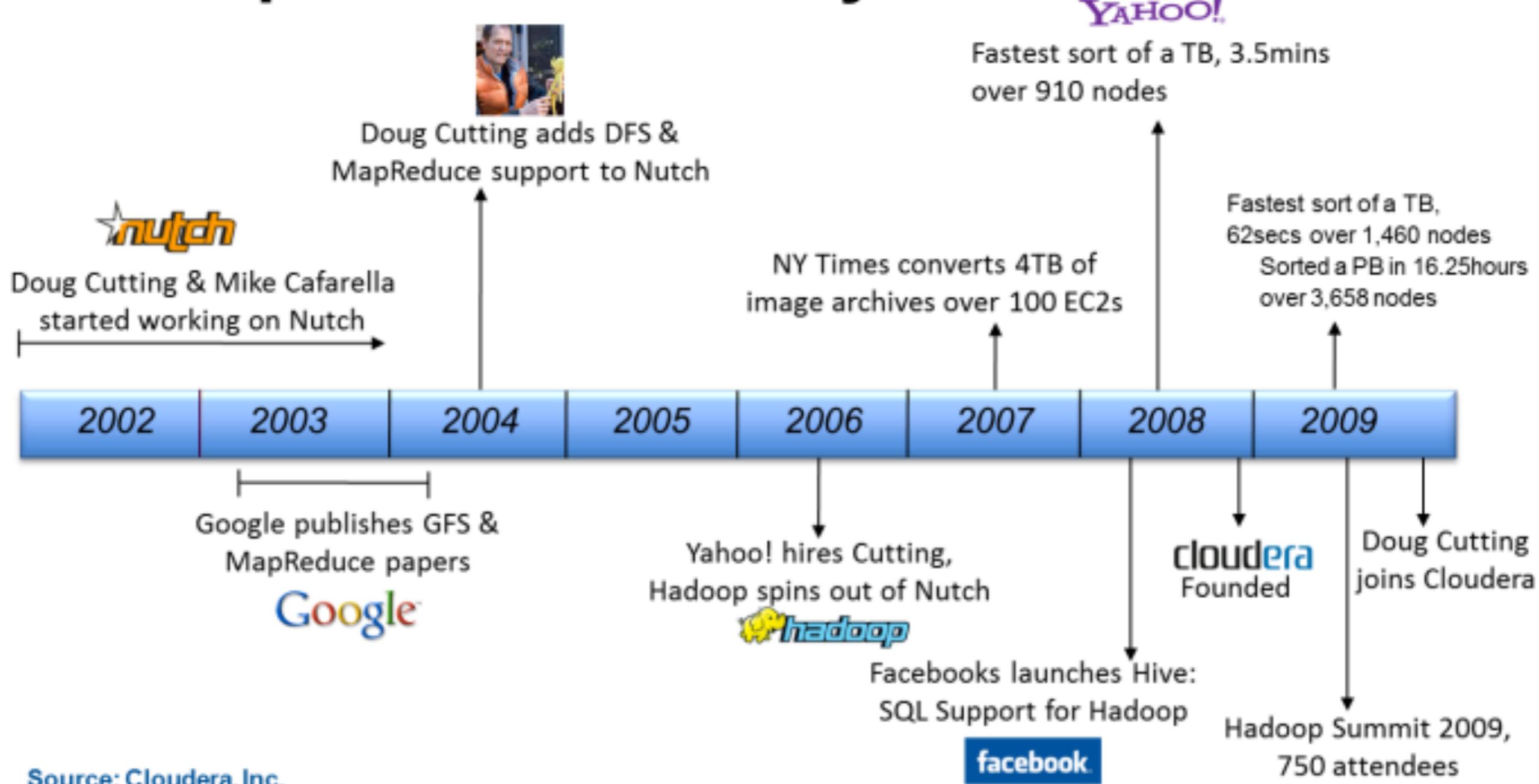
# What is Hadoop ?

A scalable fault-tolerant distributed system  
for data storage and processing

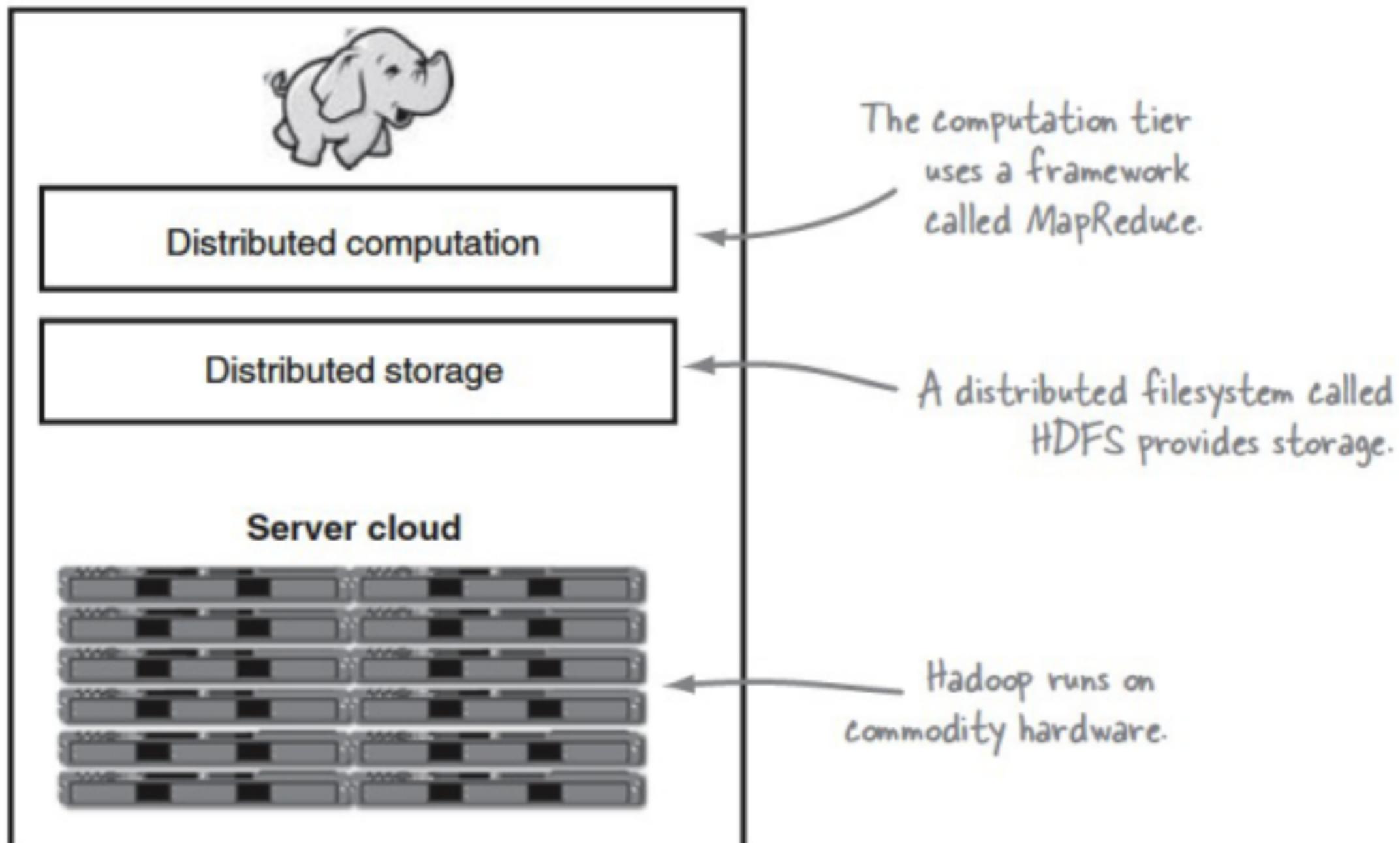
Completely written in java  
Open source & distributed under Apache license



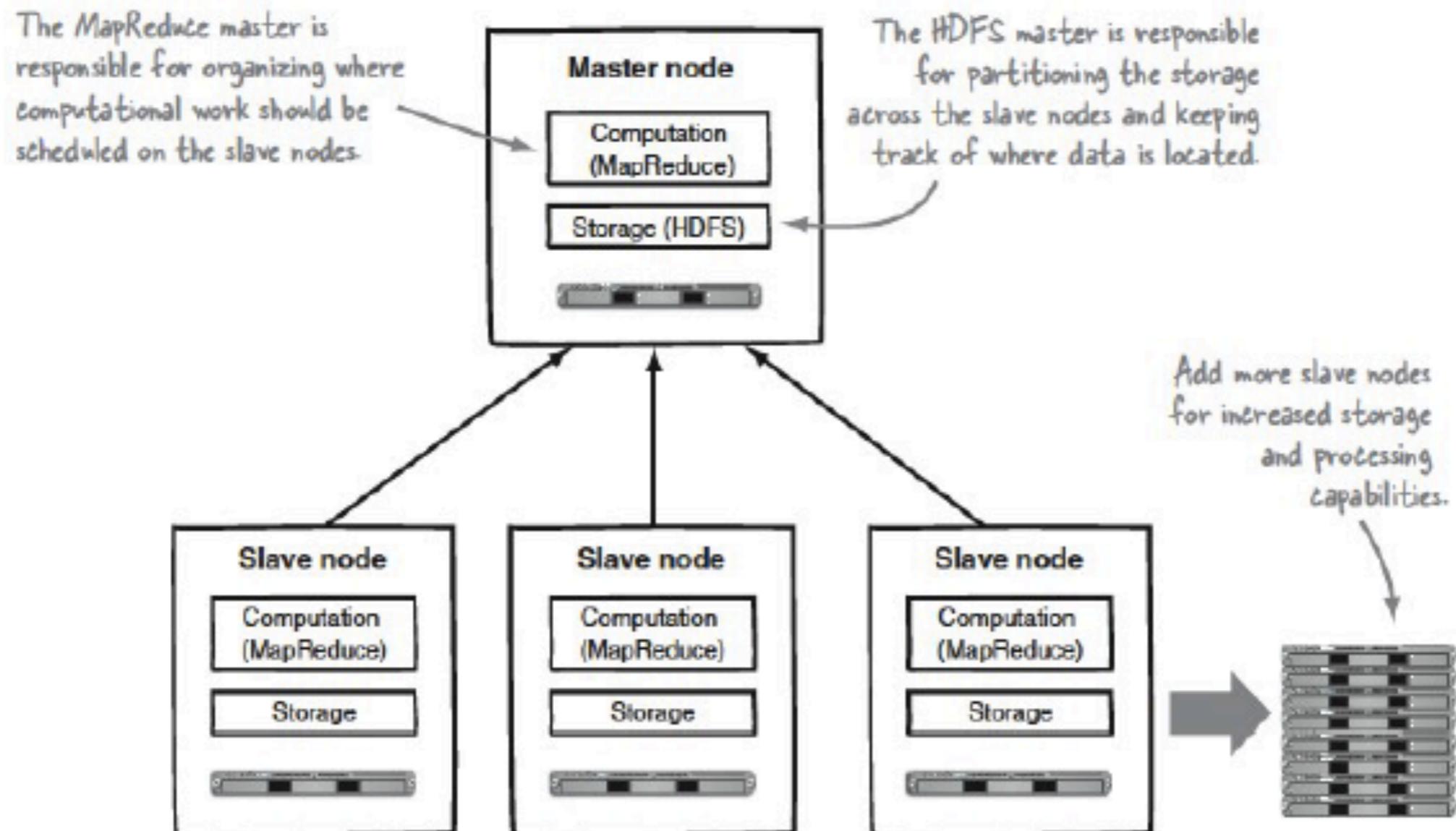
# Hadoop Creation History



# Hadoop Environment



# Hadoop Architecture



# Hadoop 2.X

## Hadoop 1

- Silos & Largely batch
- Single Processing engine



## Hadoop 2 w/YARN

- Multiple Engines, Single Data Set
- Batch, Interactive & Real-Time

**Batch**  
MapReduce

**Interactive**  
Others

**Real-Time**  
Others

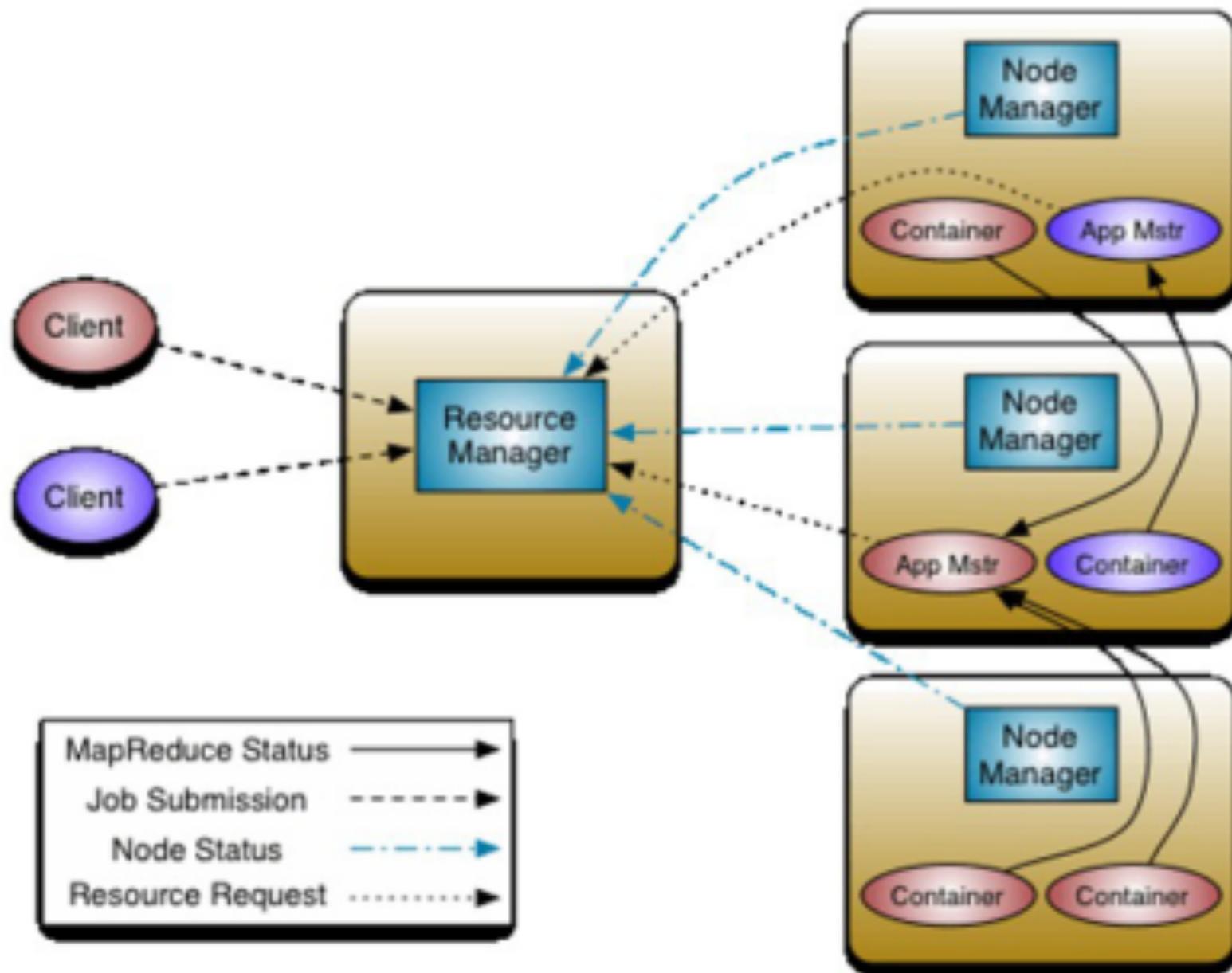
**YARN: Data Operating System**  
(Cluster Resource Management)

**MapReduce**  
(Cluster Resource Management  
& Batch Data Processing)

1 . . . . . . . .  
**HDFS**  
(Hadoop Distributed File System)

1 . . . . . . . .  
. . . . . . . . N  
**HDFS**  
(Hadoop Distributed File System)

# YARN: Yet Another Resource Negotiator



MRV2 maintains API compatibility with previous stable release (hadoop-1.x). This means that all Map-Reduce jobs should still run unchanged on top of MRv2 with just a recompile.

[Hadoop.apache.org](http://Hadoop.apache.org)

# Evolution of the Hadoop Platform

The stack is continually evolving and growing!

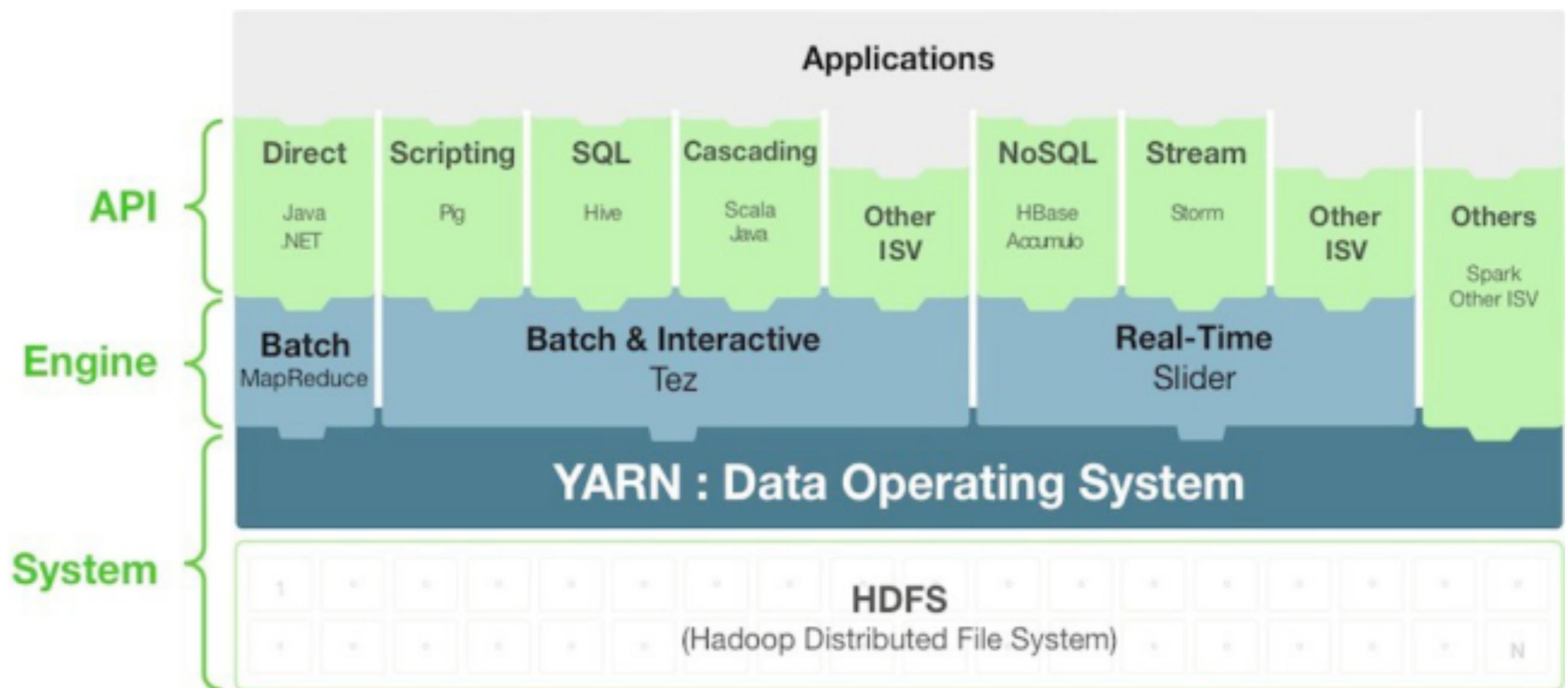
Core Hadoop (HDFS, MapReduce)	Solr Pig Core Hadoop	HBase ZooKeeper Solr Pig Core Hadoop	Hive Mahout HBase ZooKeeper Solr Pig Core Hadoop	Hive Mahout HBase ZooKeeper Solr Pig Core Hadoop	Mahout HBase ZooKeeper Solr Pig YARN Core Hadoop	Flume Bigtop Oozie MRUnit HCatalog Hue Sqoop Whirr Avro Hive Mahout HBase ZooKeeper Solr Pig YARN Core Hadoop	Spark Tez Impala Kafka Drill Flume Bigtop Oozie MRUnit HCatalog Hue Sqoop Whirr Avro Hive Mahout HBase ZooKeeper Solr Pig YARN Core Hadoop	Parquet Sentry Spark Tez Impala Kafka Drill Flume Bigtop Oozie MRUnit HCatalog Hue Sqoop Whirr Avro Hive Mahout HBase ZooKeeper Solr Pig YARN Core Hadoop	Ibis Flink Parquet Sentry Spark Tez Impala Kafka Drill Flume Bigtop Oozie MRUnit HCatalog Hue Sqoop Whirr Avro Hive Mahout HBase ZooKeeper Solr Pig YARN Core Hadoop
2006	2007	2008	2009	2010	2011	2012	2013	2014-15	

cloudera

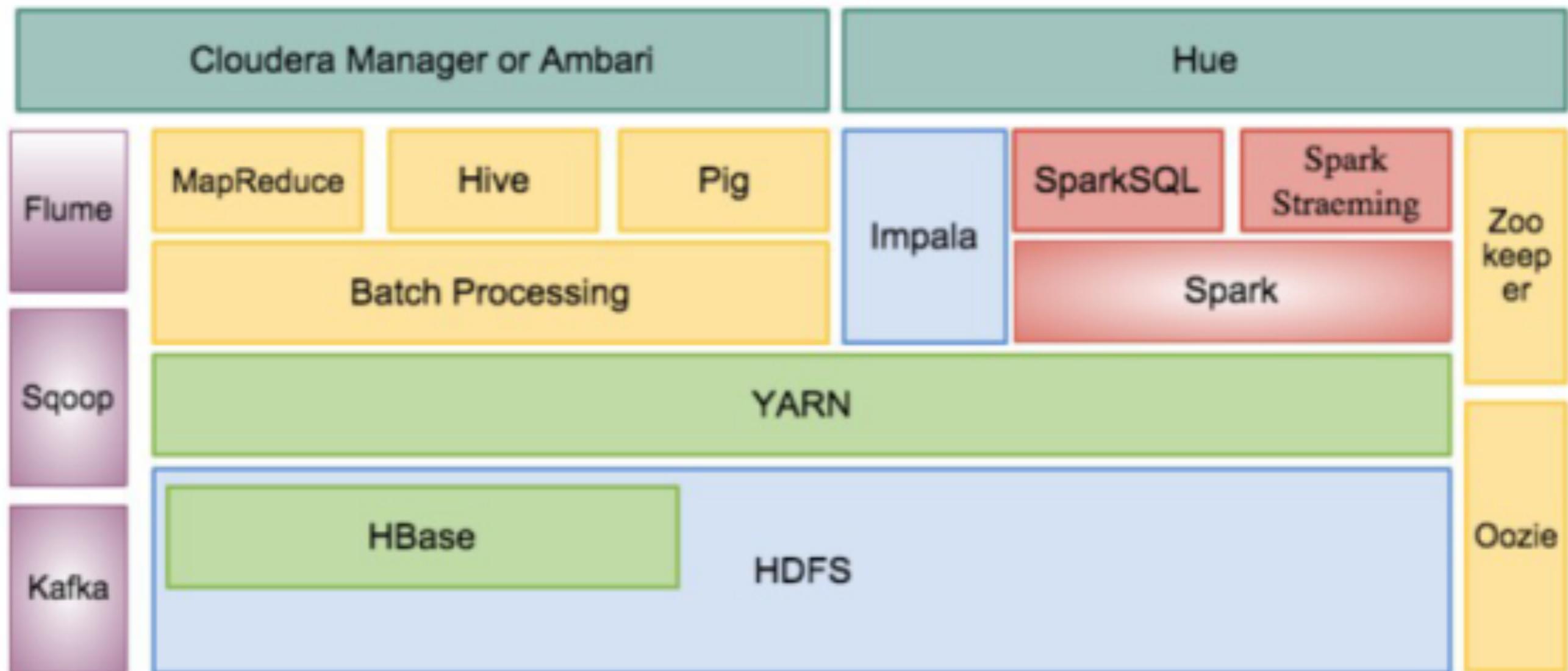
© Cloudera, Inc. All rights reserved.

9

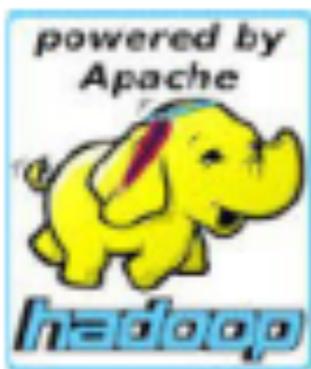
# Hadoop 2.x Ecosystems



# Hadoop Ecosystems



# Hadoop Distribution



MAPR

cloudera



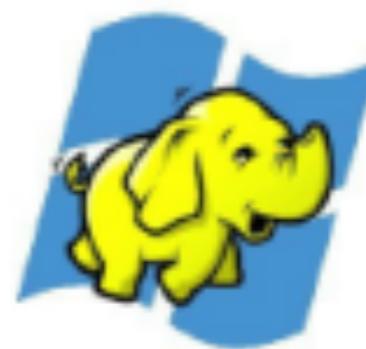
Pivotal™



TERADATA.



amazon  
web services™

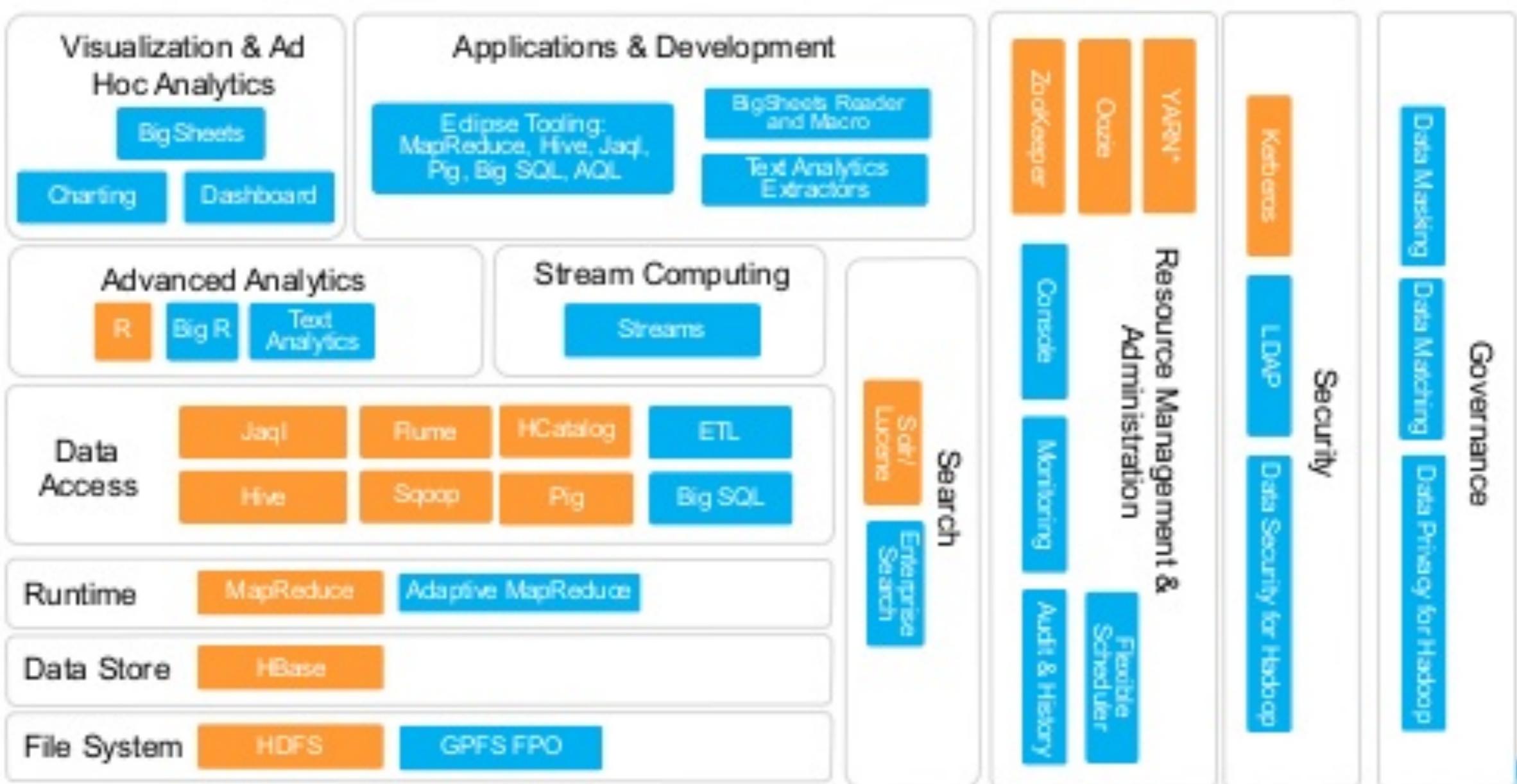


Microsoft Azure

# IBM InfoSphere BigInsights



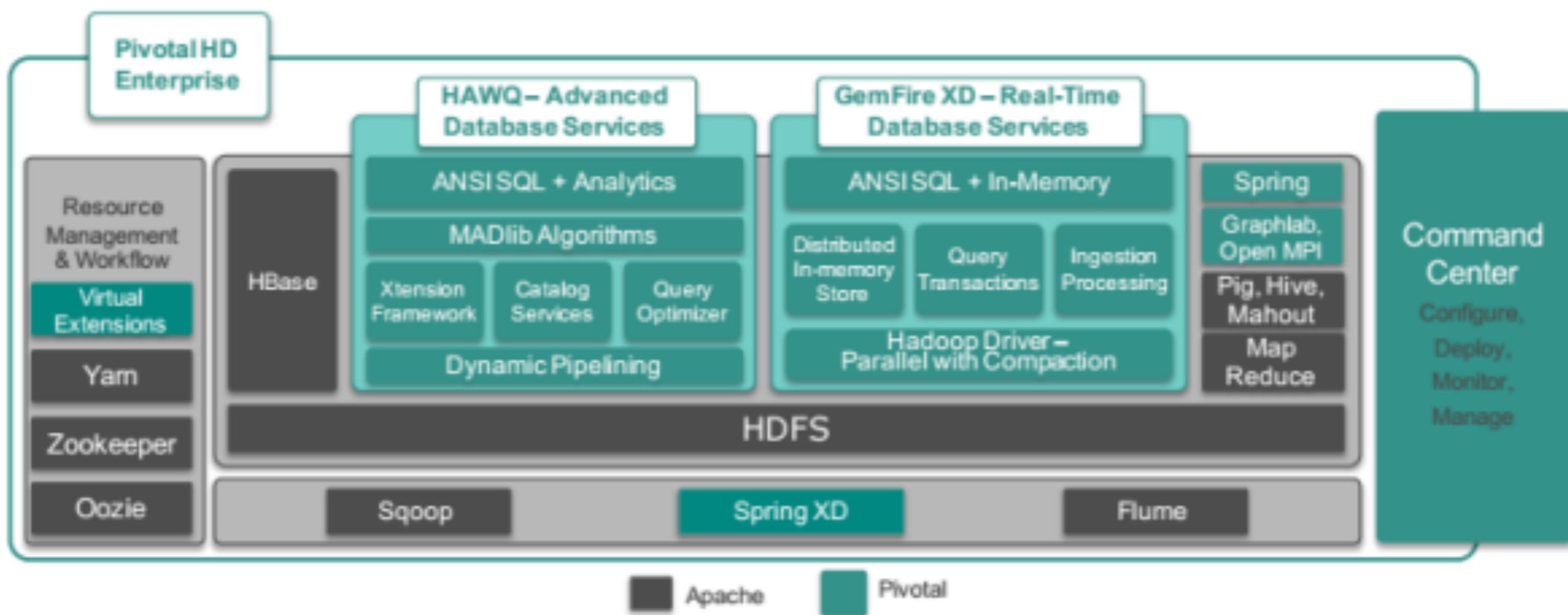
## IBM InfoSphere BigInsights for Hadoop



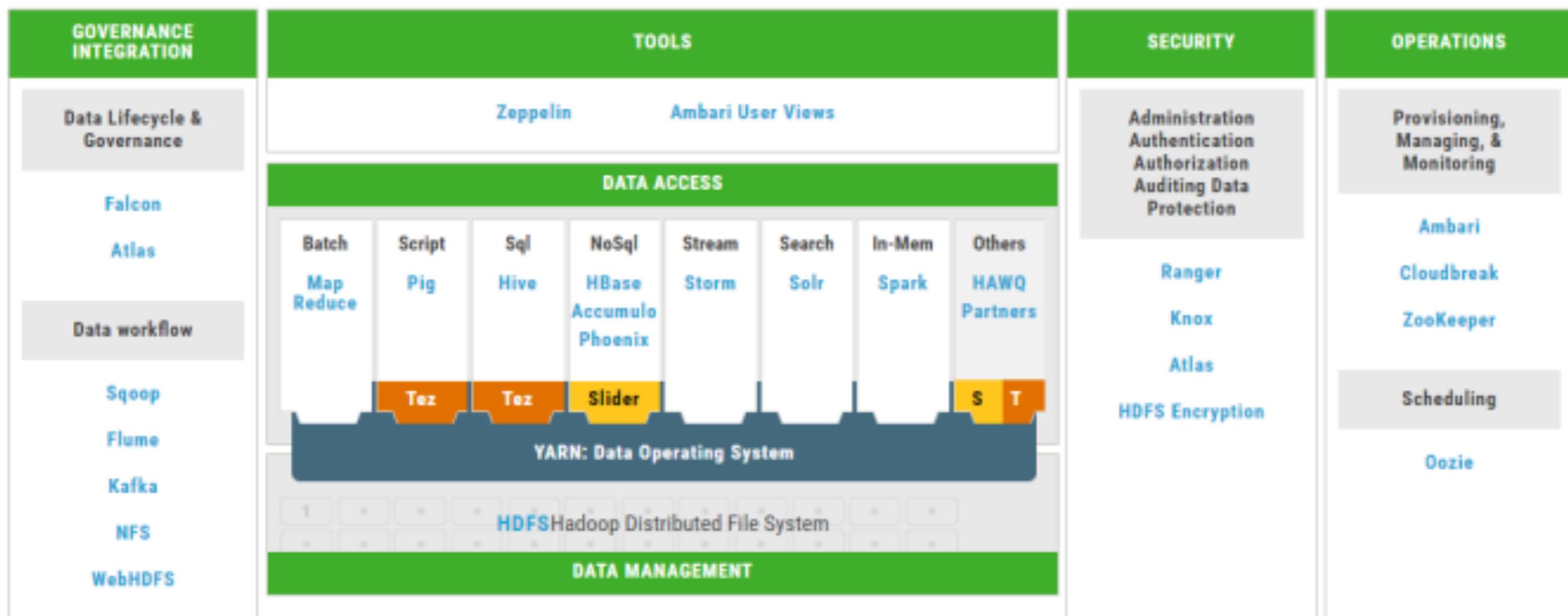
\* In Beta

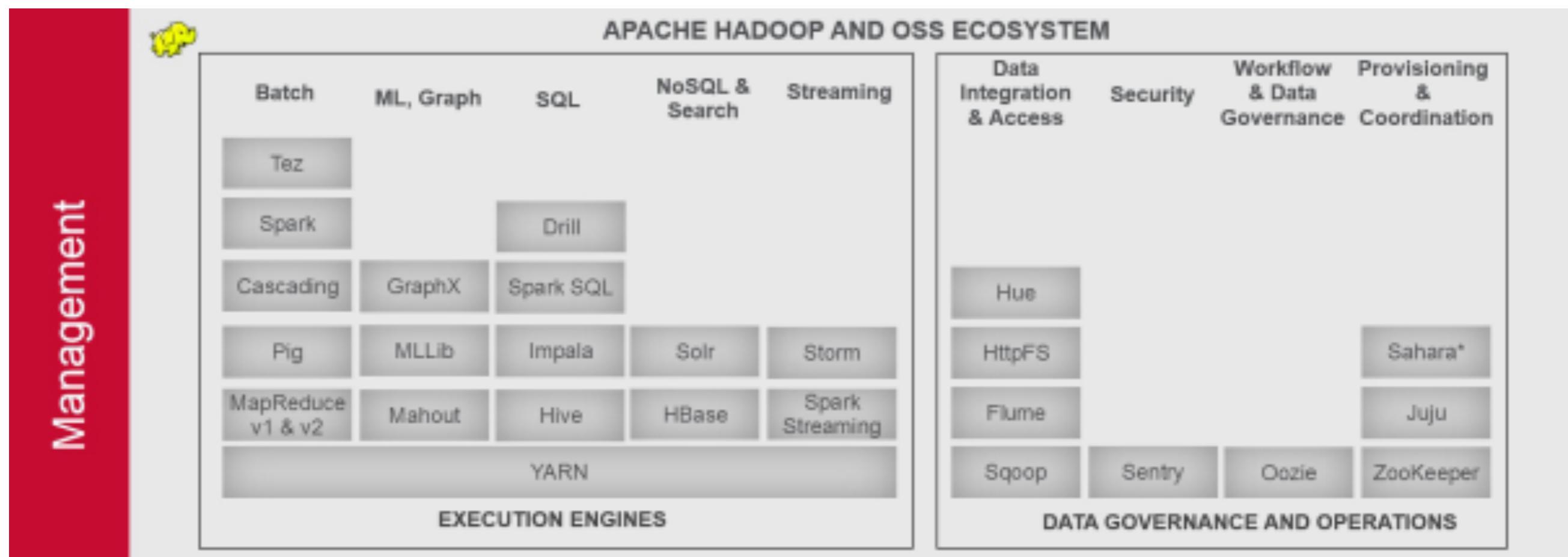


# Pivotal HD Architecture



# Hortonworks

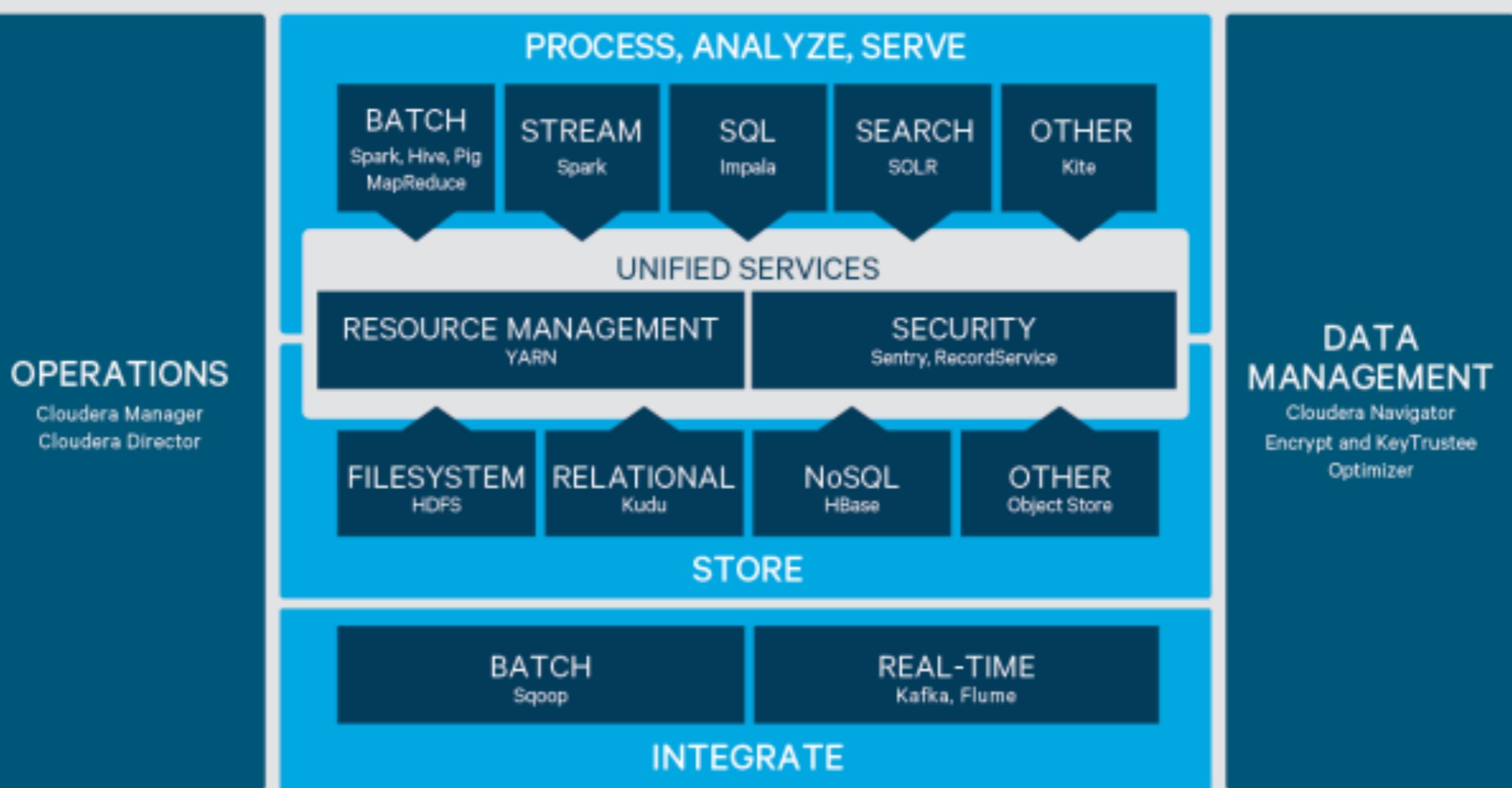




MapR-FS

Data Platform

MapR-DB



# Default Cloudera Services

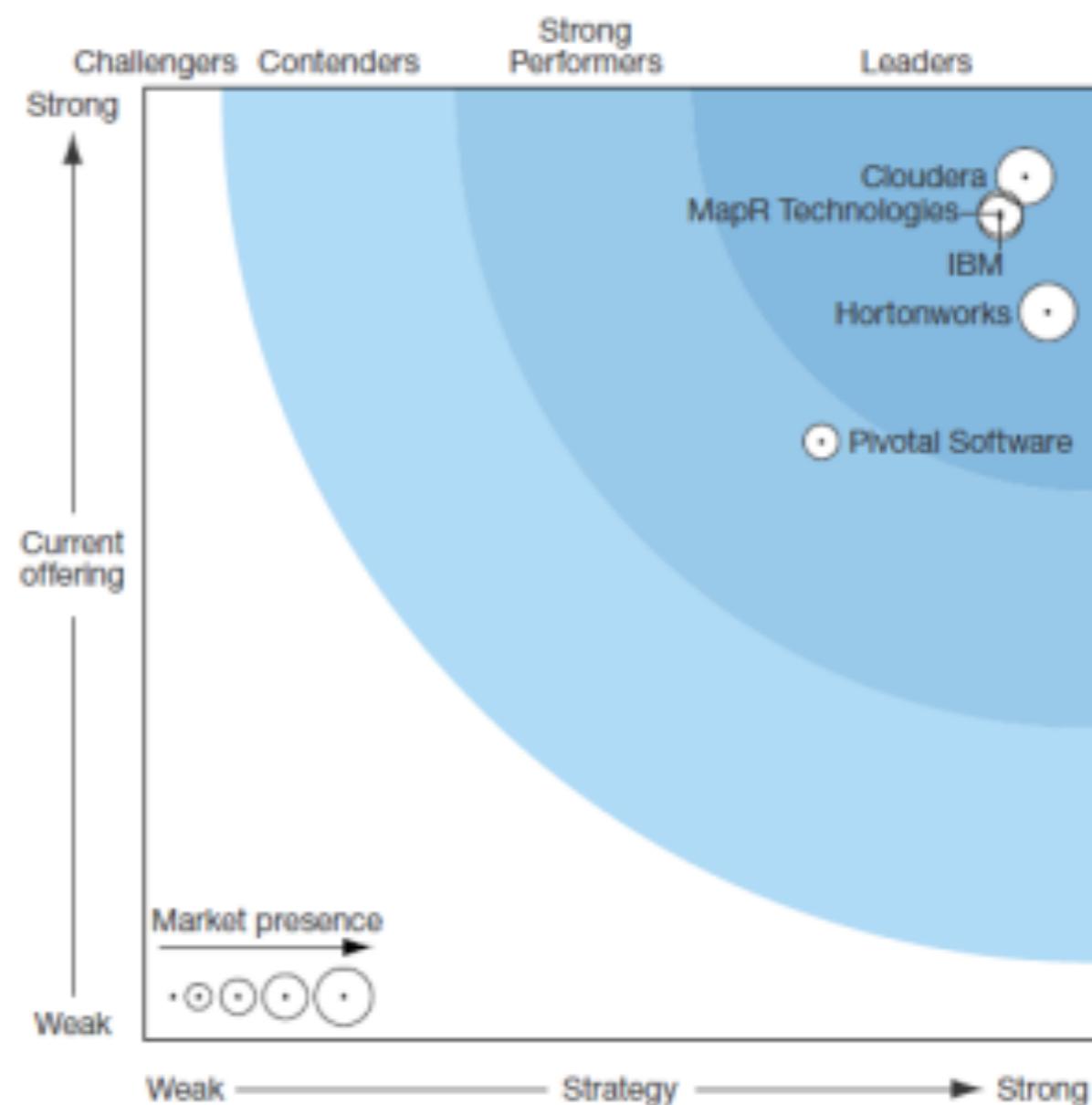
- Cloudera Manager
- HDFS
- YARN
- Apache Hive and Pig
- Apache Flume and Sqoop
- Apache Oozie
- Cloudera Hue
- ZooKeeper



**cloudera**

## Big Data Hadoop Distributions Q1 2016

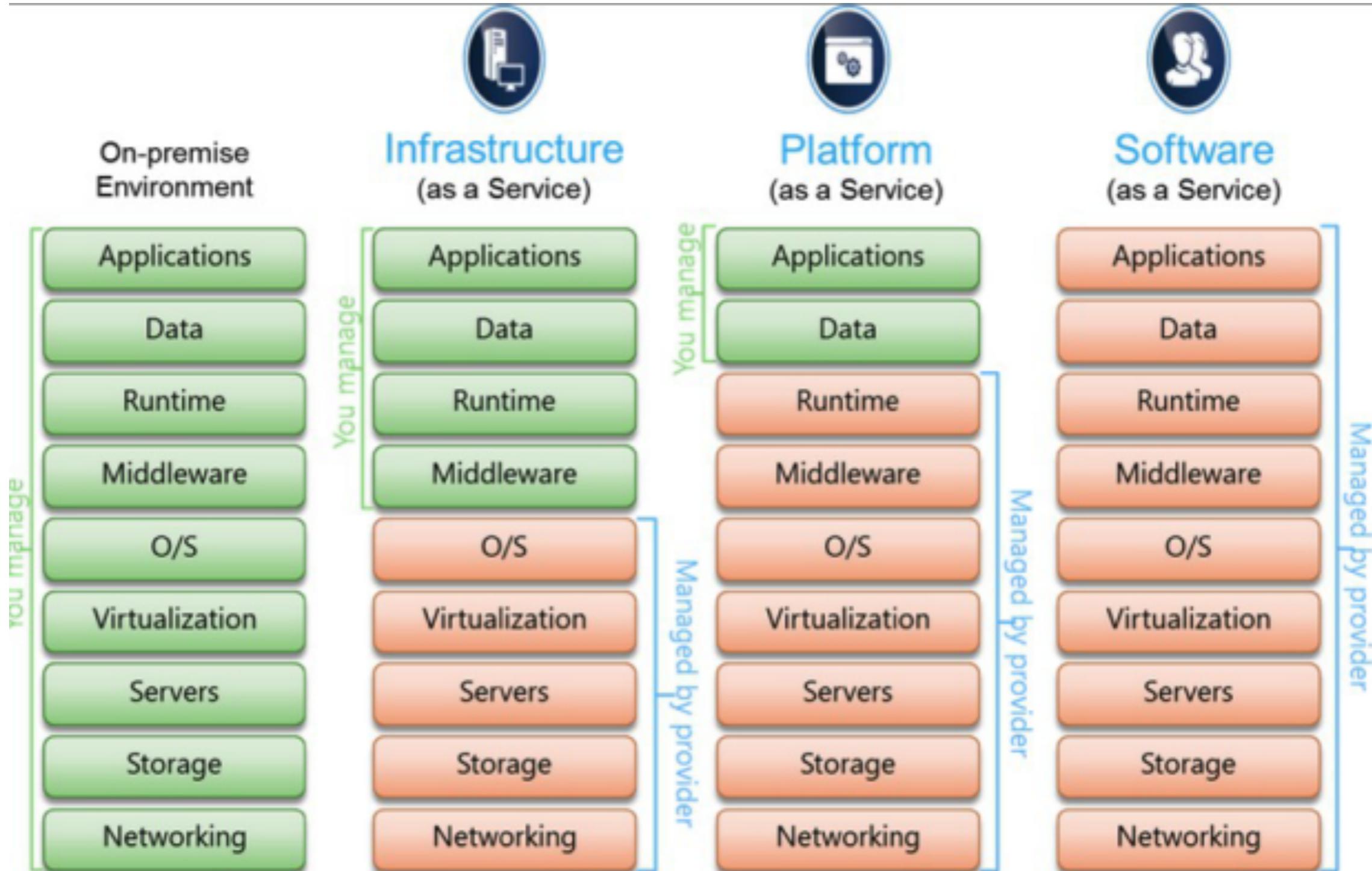
Vendor	Product evaluated	Product version evaluated
Cloudera	Cloudera Enterprise	5.50
Hortonworks	Hortonworks Data Platform	2.30
IBM	IBM BigInsights for Apache Hadoop	4.10
MapR Technologies	The MapR Distribution including Apache	5.00
Pivotal Software	HadoopPivotal HD	3.x

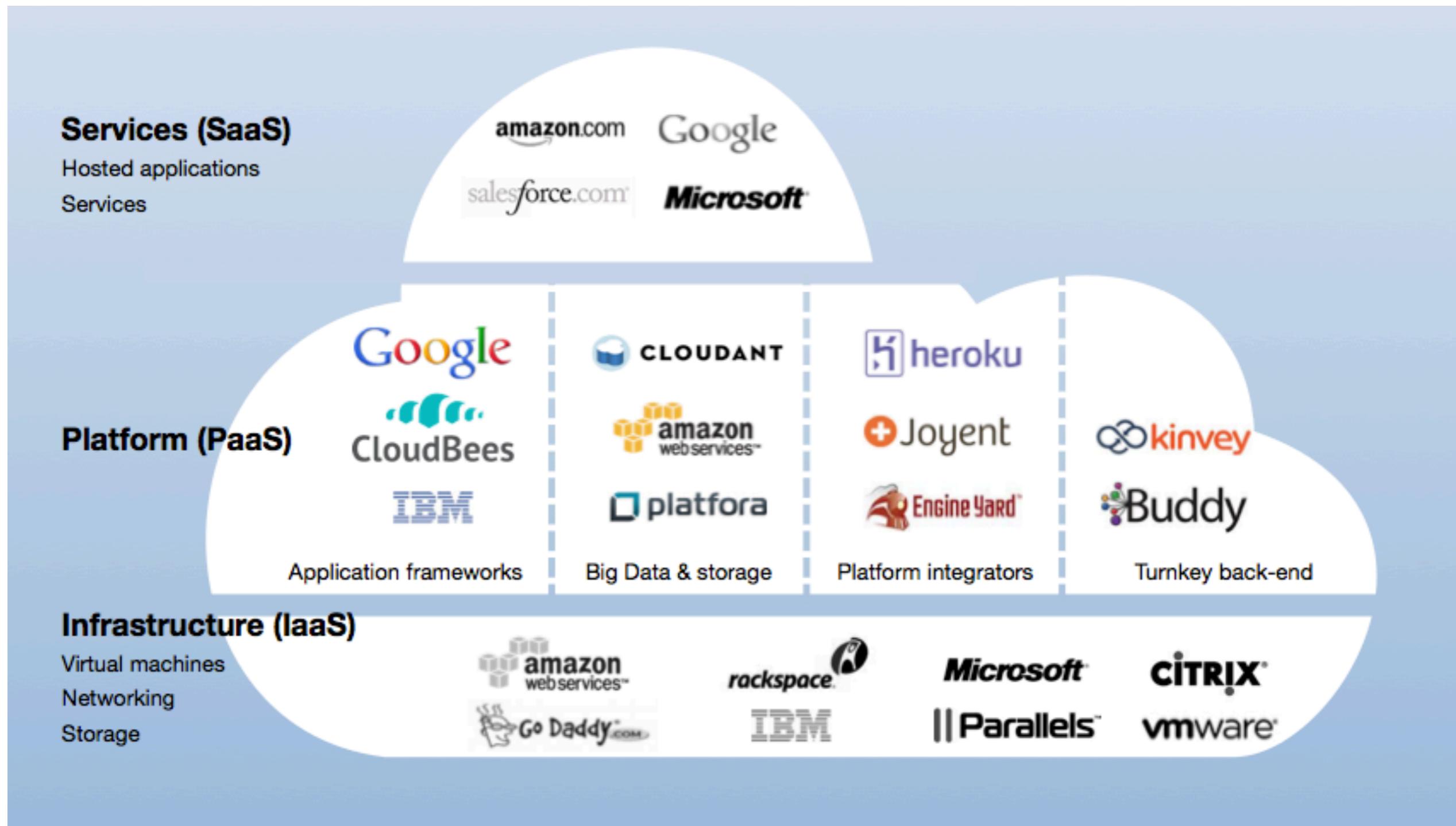


## Issue with Big Data Infrastructure

- Large investment
- Scalability
- ROI
- Business Cases

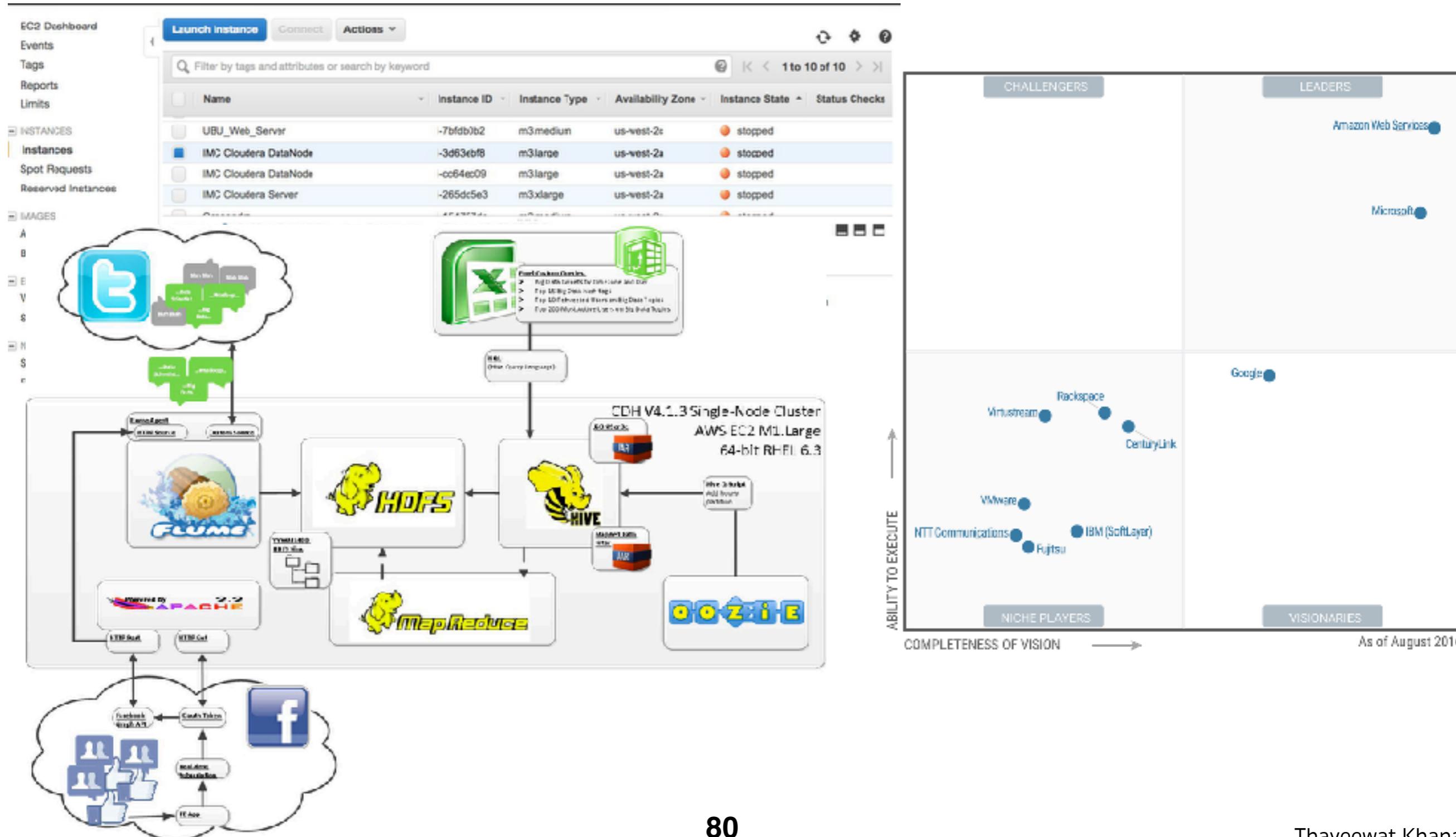
# Cloud Technology





# Big Data on Cloud

## Big Data using IaaS

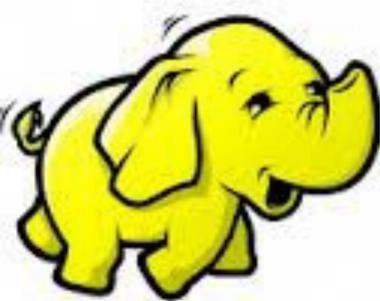


# Big Data on Cloud

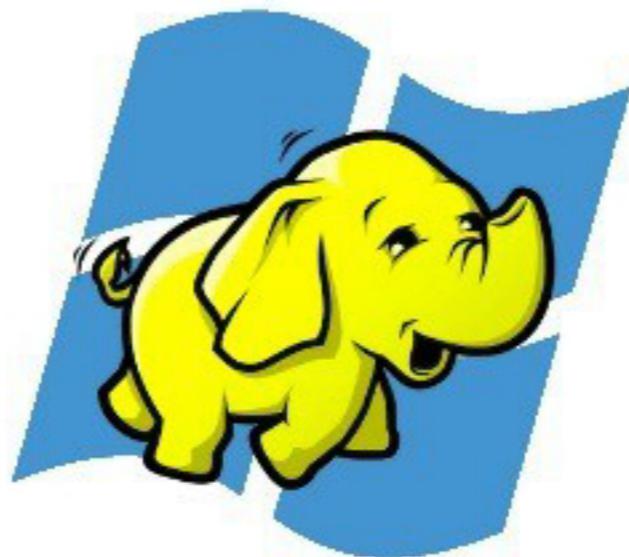
## Using Big Data as a Services



Amazon



Elastic Mapreduce



Microsoft Azure Hadoop