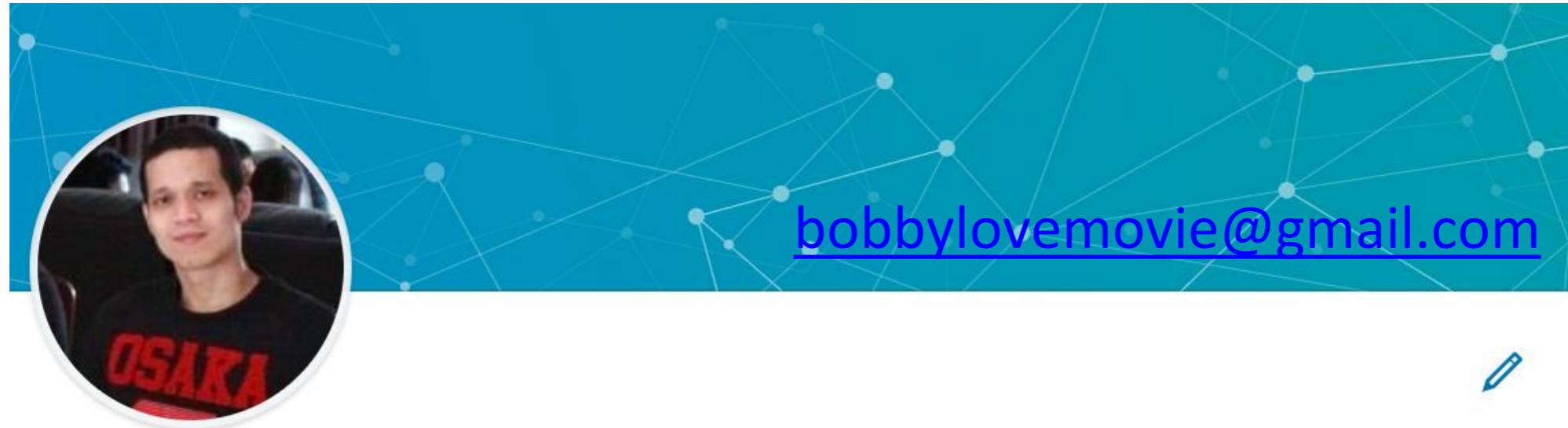


Introduction To Big Data

<http://bit.ly/2v7OZTh>



Bobby Thaveewat



Dhipaya Insurance PCL (TIP)



Data Science Specialist

Dhipaya Insurance PCL (TIP)

Jul 2018 – Present · 10 mos

Bangkok Metropolitan Area, Thailand



Senior System Analyst(IT-CIS Finance Delivery Big Data)

True Corporation

Jun 2016 – Present · 2 yrs 11 mos

Bangkok Metropolitan Area, Thailand



Programmer Analysis, System Analysis

CDG Systems (CDGS)

May 2005 – Present · 14 yrs

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



"นายกฯ" ชี้ ใช้ "Big Data" บริหารราชการแผ่นดิน

15 May 2018

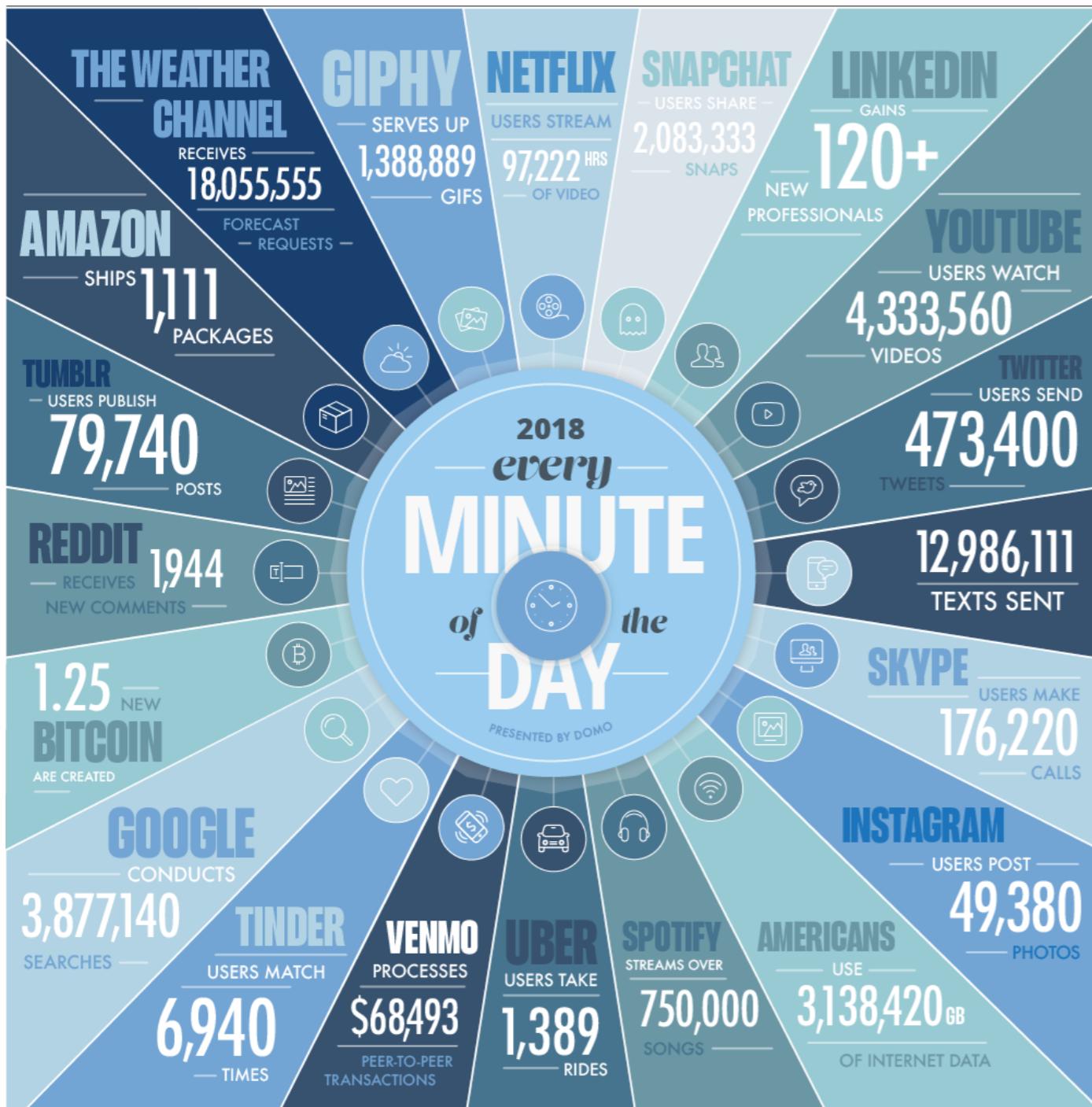
What is Big Data ?

“ 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.

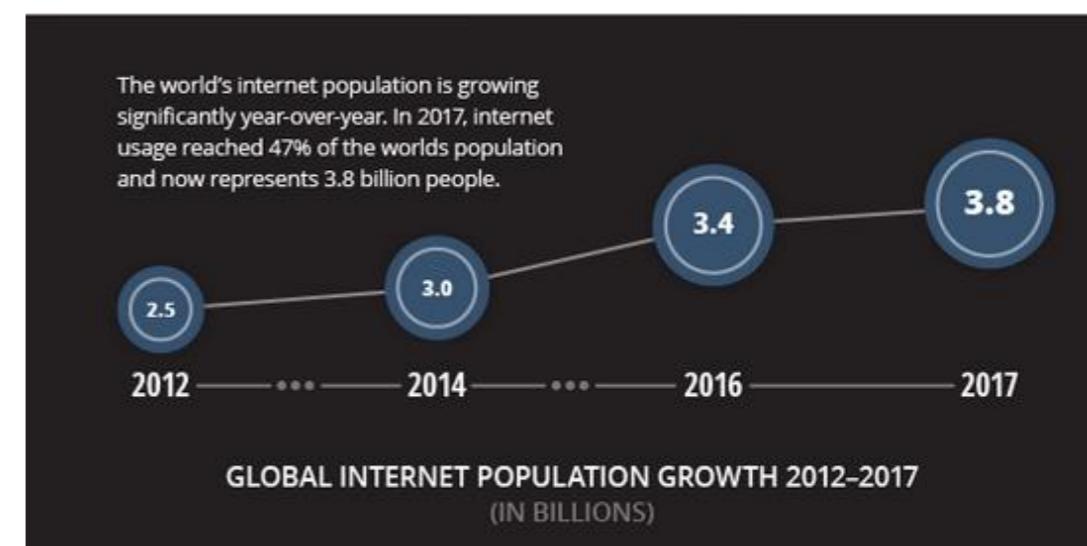
To gain value from this data,
you must choose an alternative way to process it. ”

Big Data Now: O'Reilly Media

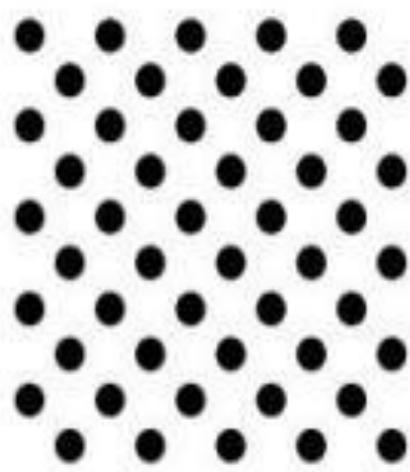
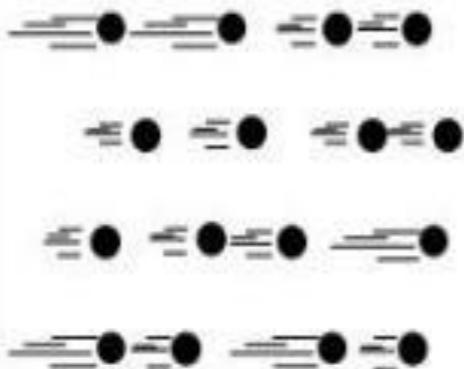
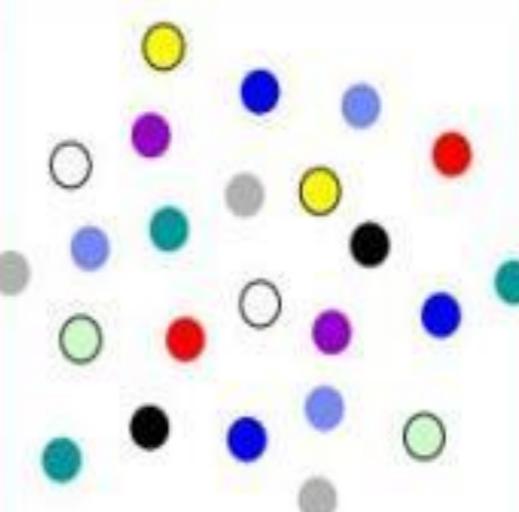
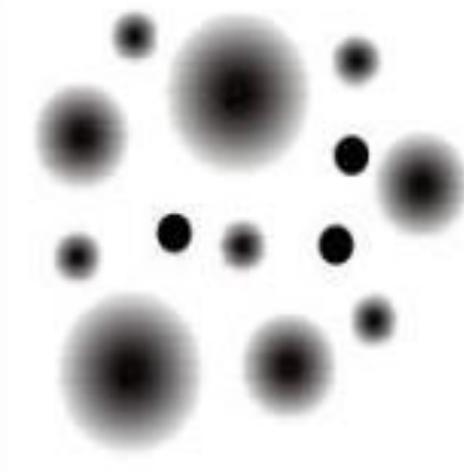
Why Big Data ?



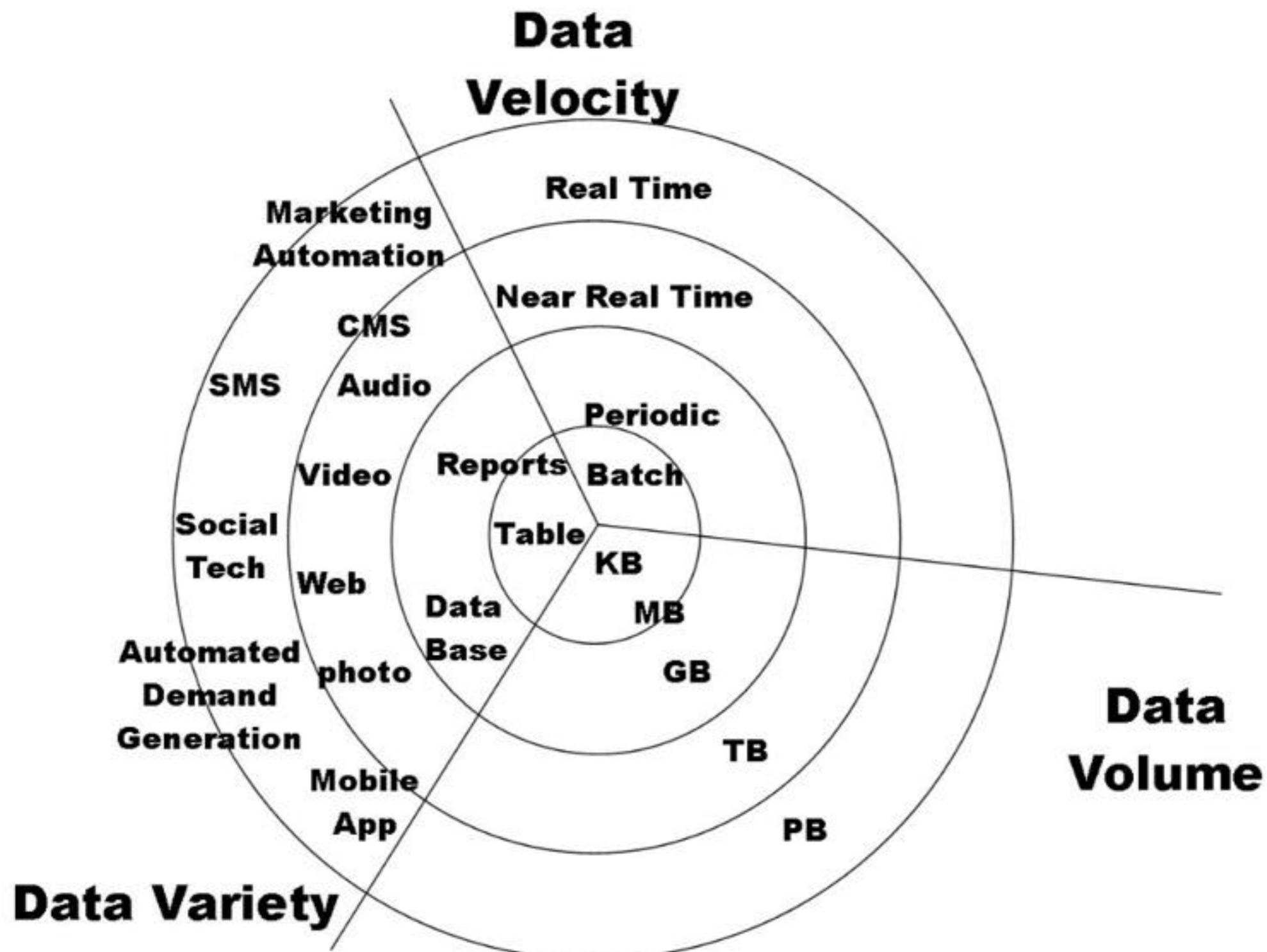
DATA NEVER SLEEP 6.0

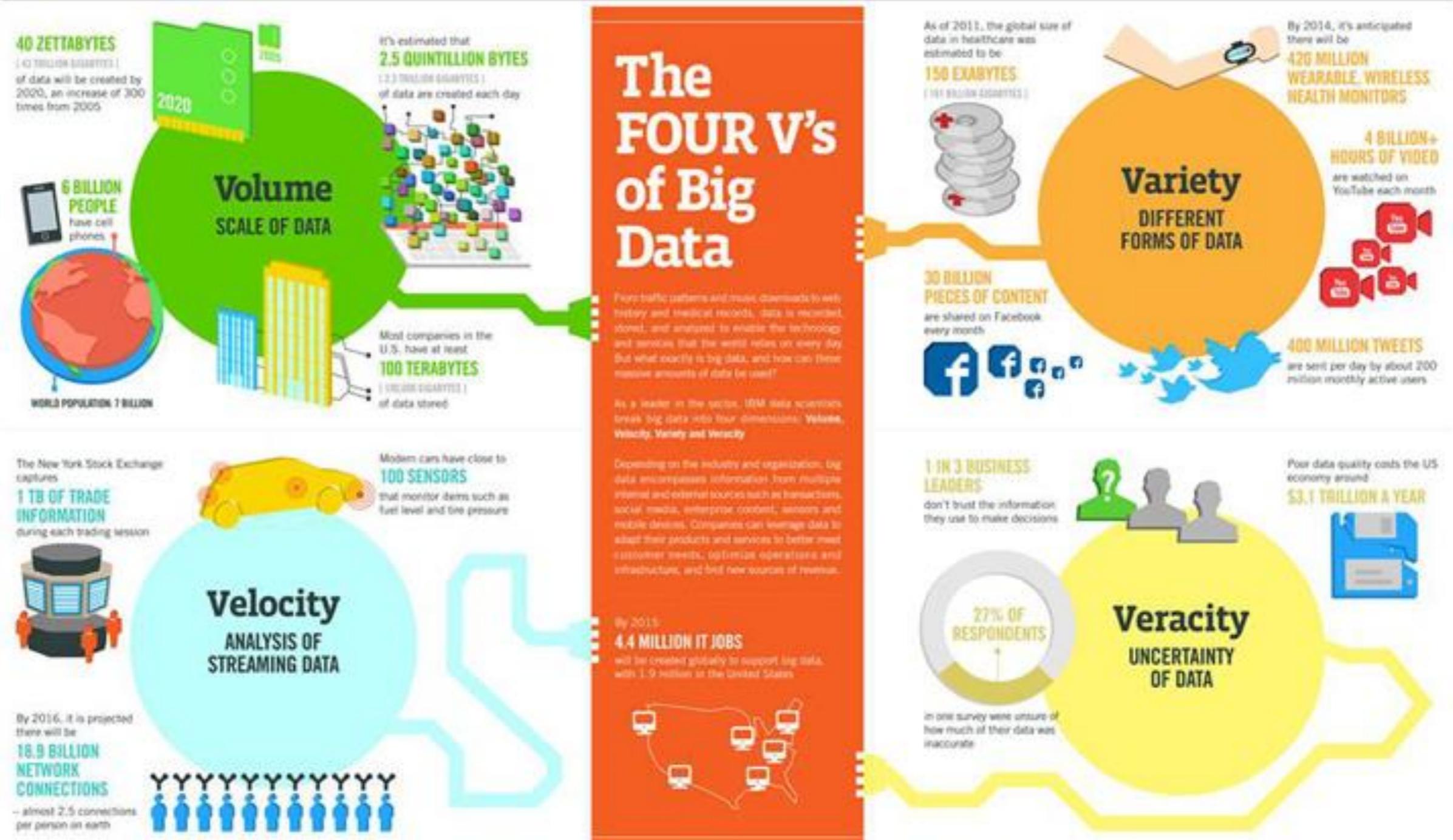


Three Characteristics of Big Data

Volume	Velocity	Variety	Veracity*
			
Data at Rest	Data in Motion	Data in Many Forms	Data in Doubt
Terabytes to exabytes of existing data to process	Streaming data, milliseconds to seconds to respond	Structured, unstructured, text, multimedia	Uncertainty due to data inconsistency & incompleteness, ambiguities, latency, deception, model approximations

Source Introduction to Big Data: Dr. Putchong Uthayopas





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



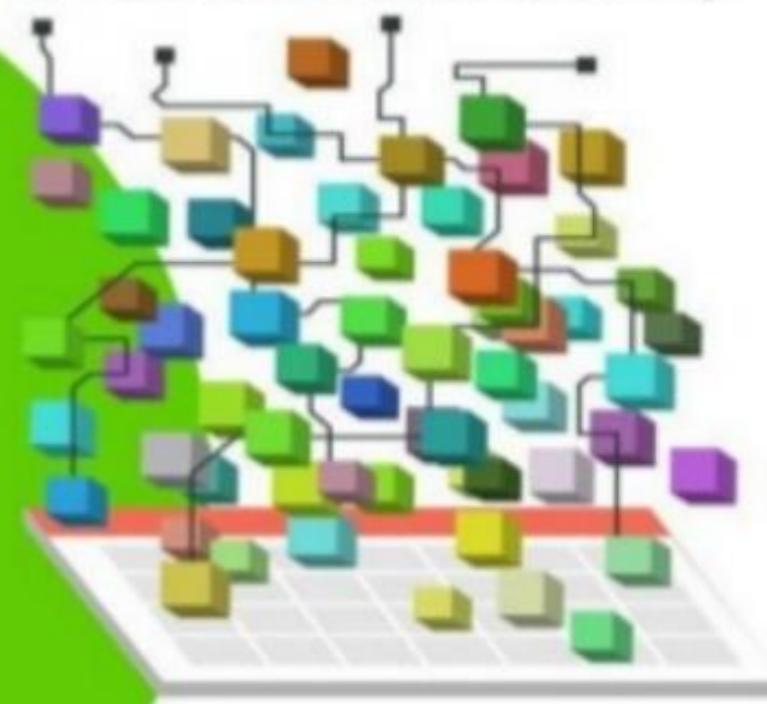
WORLD POPULATION: 7 BILLION



Volume SCALE OF DATA



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

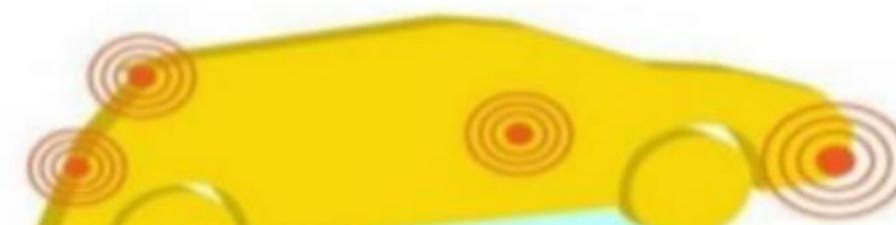
during each trading session



By 2016, it is projected there will be

18.9 BILLION NETWORK CONNECTIONS

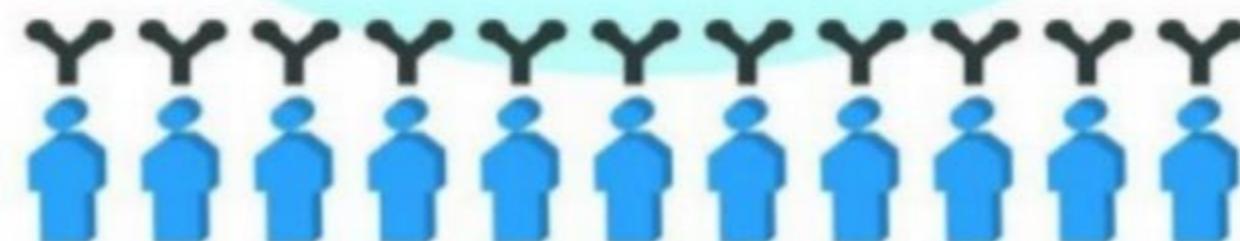
– almost 2.5 connections per person on earth



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



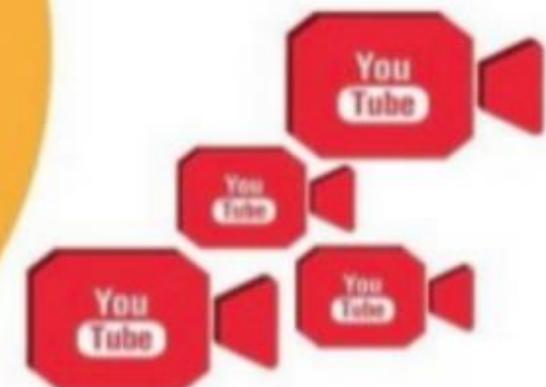
By 2014, it's anticipated there will be

**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

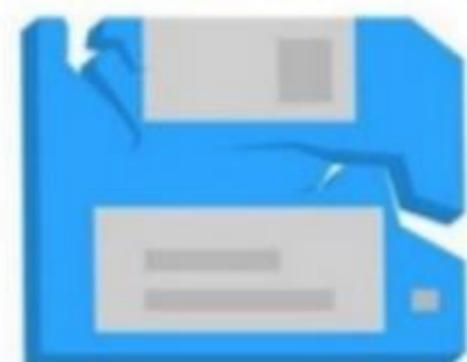


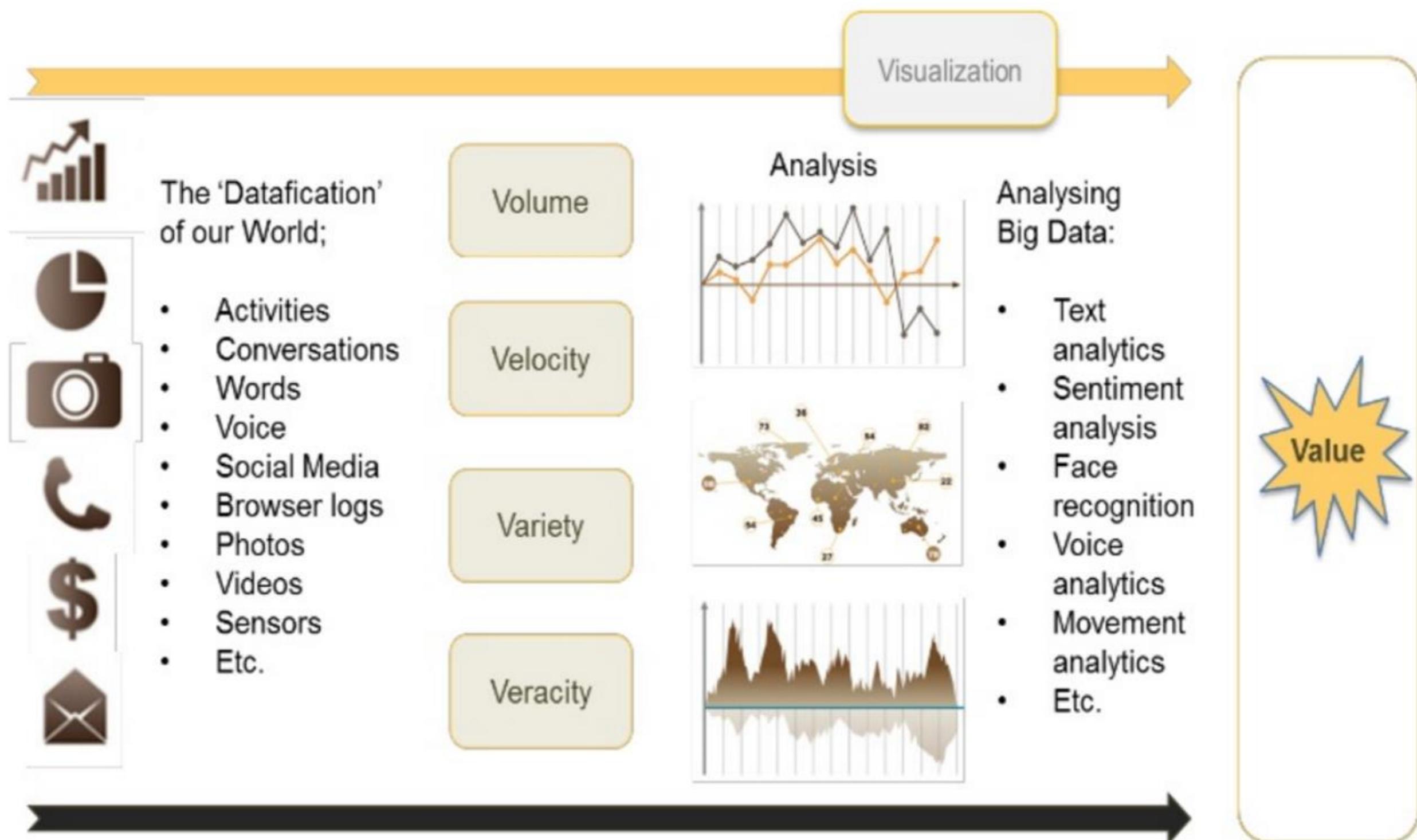
27% OF
RESPONDENTS

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

Veracity UNCERTAINTY OF DATA

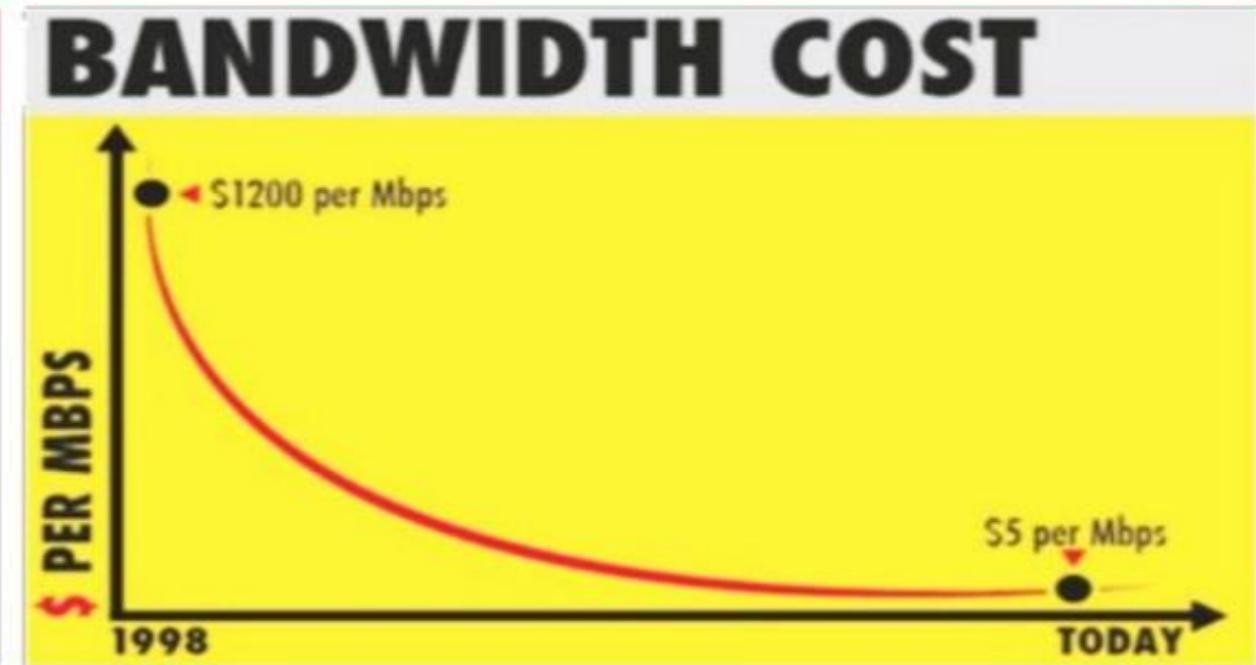
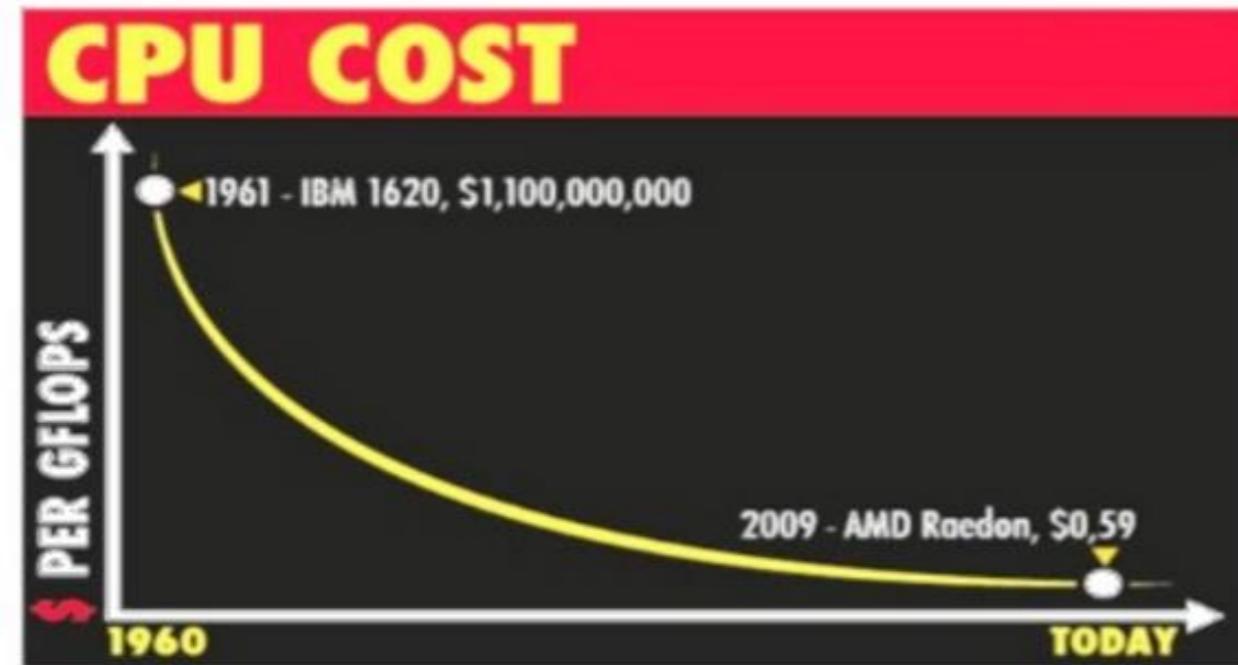
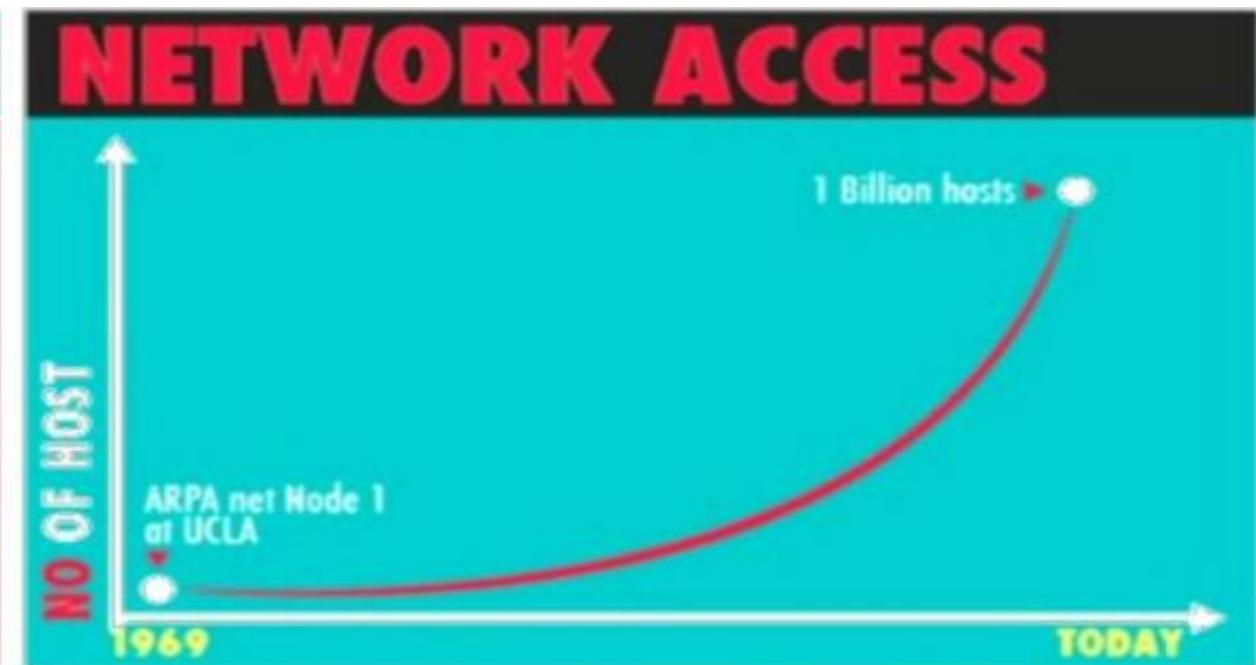
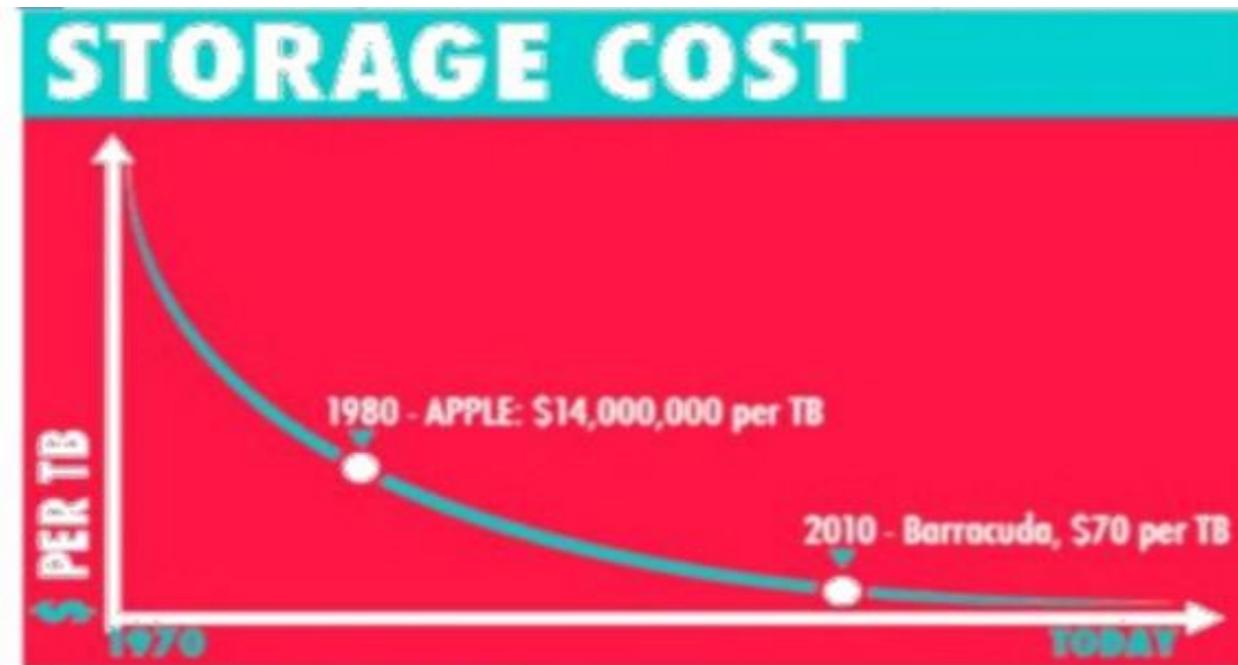
Poor data quality costs the US
economy around
\$3.1 TRILLION A YEAR





Source: Bernard Marr

Big Data : Why Now?

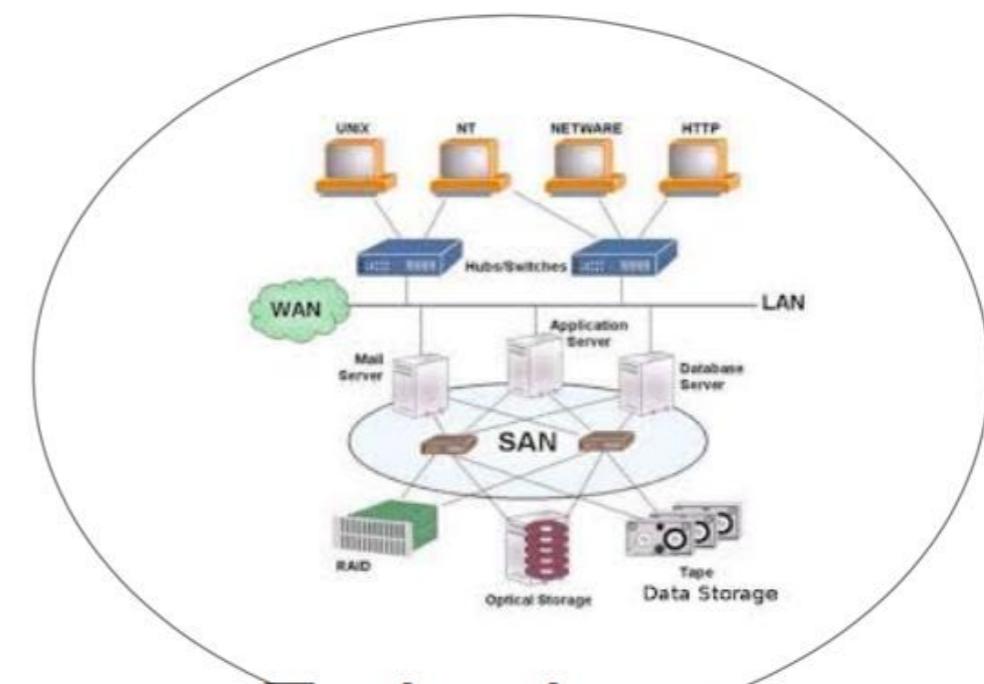


ตัวอย่างอุตสาหกรรมที่มีข้อมูลขนาดใหญ่

- **โทรคมนาคม** ให้บริการโทรศัพท์เคลื่อนที่ก็อาจมีข้อมูลลูกค้า หรือข้อมูลการใช้งานโทรศัพท์ของลูกค้า (CDR: Call Detail Record) ที่มากกว่า 1 ล้านวันอาจเป็นหลาย Terabyte หรือหลายหมื่นล้านเรคอร์ด
- **การเงินการธนาคาร** มีข้อมูลการทำธุรกรรมของลูกค้าที่มีจำนวนมาก หรือข้อมูลการซื้อขายหลักทรัพย์ในตลาดทุน
- **ค้าปลีก** ทางสังคมที่มีรายการซื้อสินค้า ข้อมูลสต็อกสินค้า ข้อมูลลูกค้า ที่มีการเคลื่อนไหวในแต่ละวันจำนวนมาก
- **การแพทย์** ข้อมูลการรักษาพยาบาลของผู้ป่วย ข้อมูลด้านยา ข้อมูลโรคระบาดต่างๆ



Data Sources



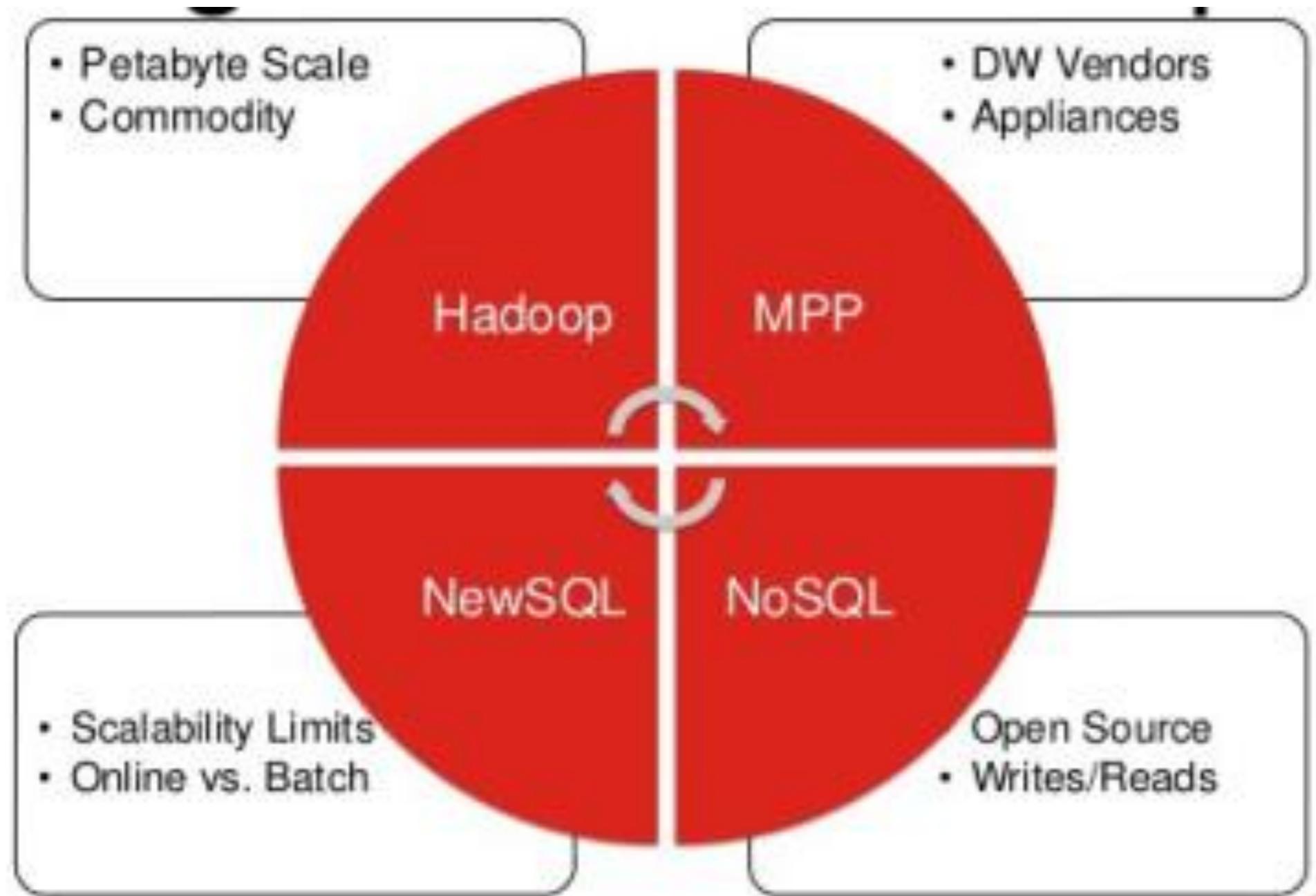
Technology



Analytics

Big Data Technology

Big Data Landscape

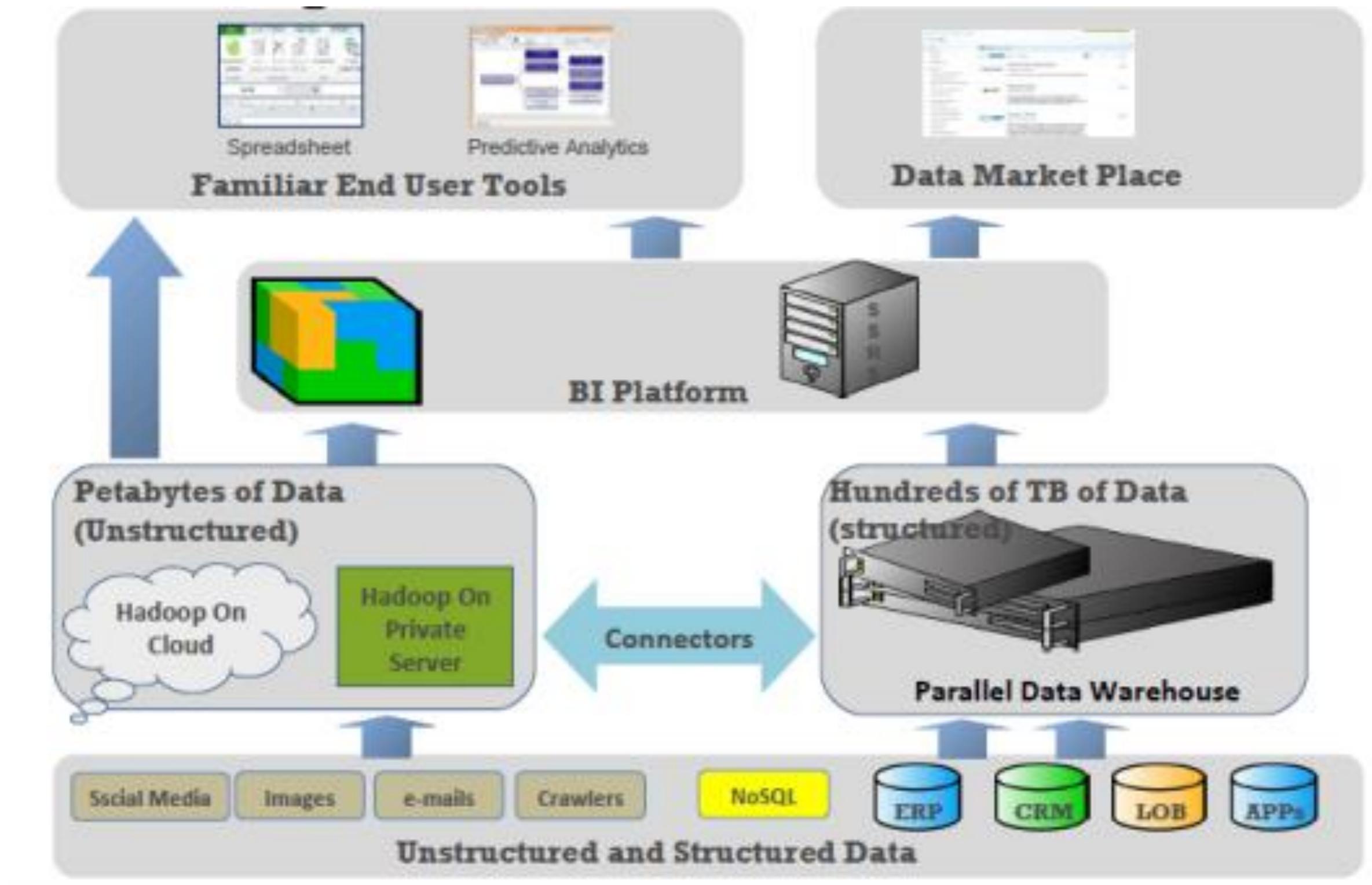




BIG DATA & AI LANDSCAPE 2018



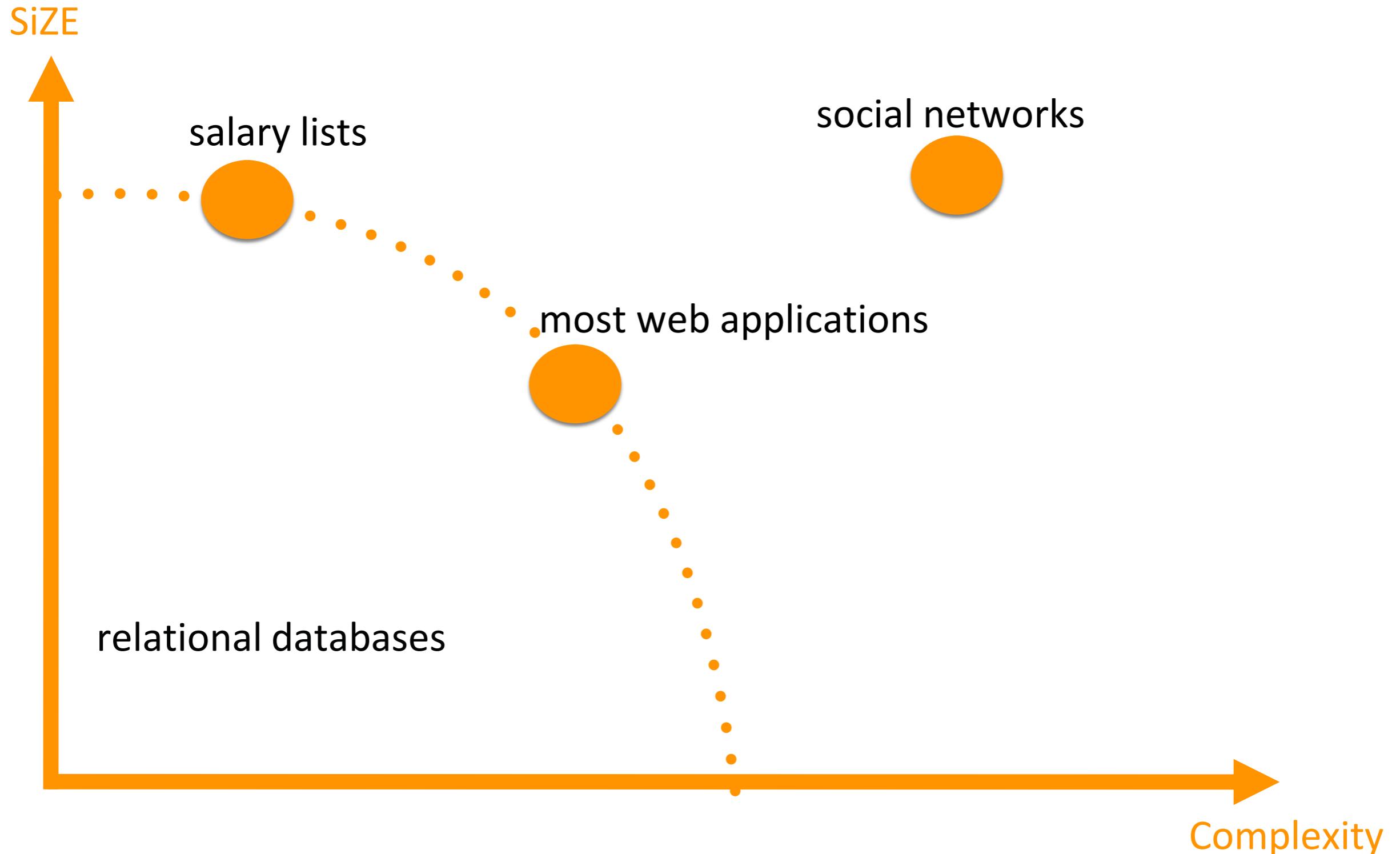
Big Data Future Architecture



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

Key Value Store

Column-oriented

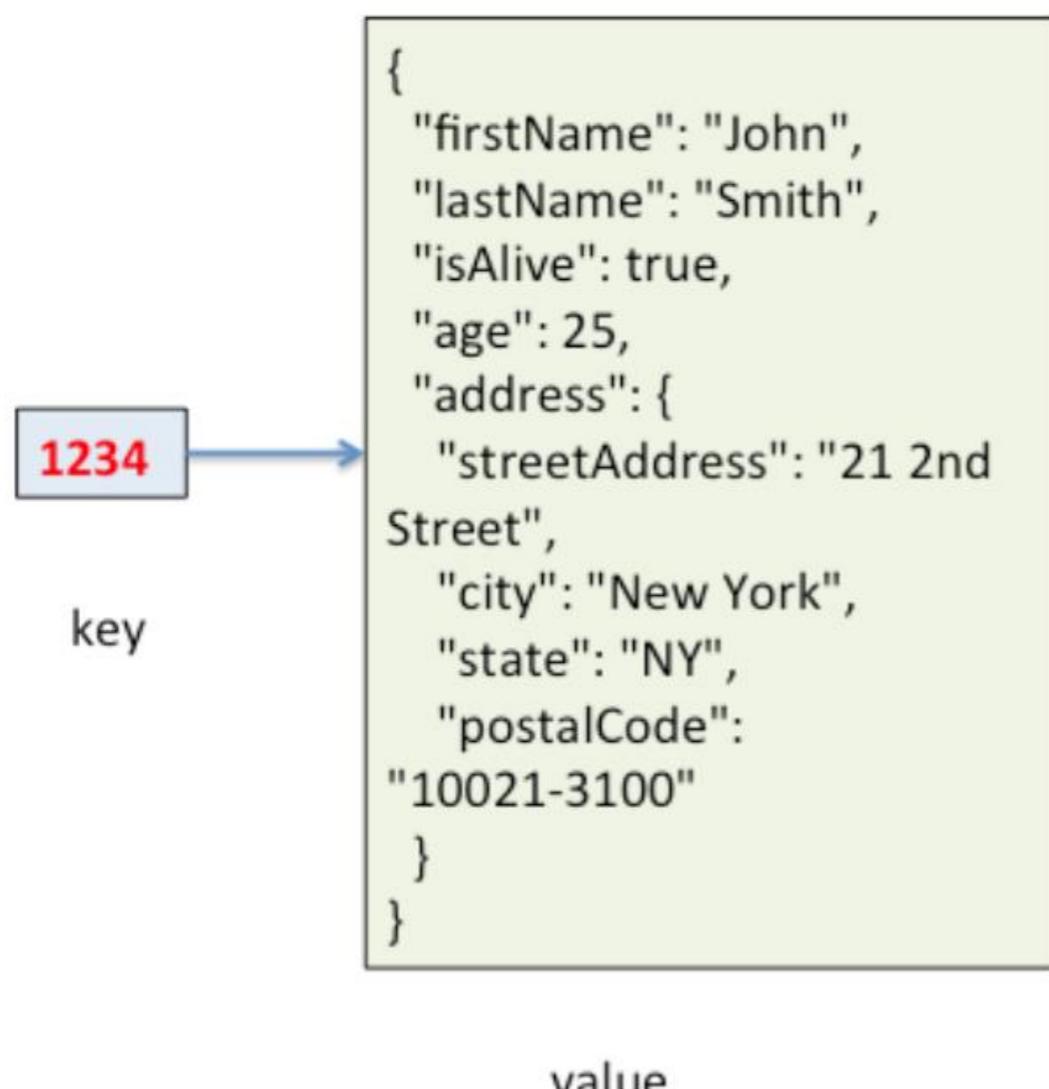
Document Store

Graph

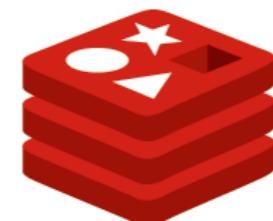
Key-value store database

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



Column-oriented databases

Row Oriented
(RDBMS Model)

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

Multi-valued

null

Column Oriented
(Multi-value sorted map)

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



Document-oriented database

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

Most of the databases available 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

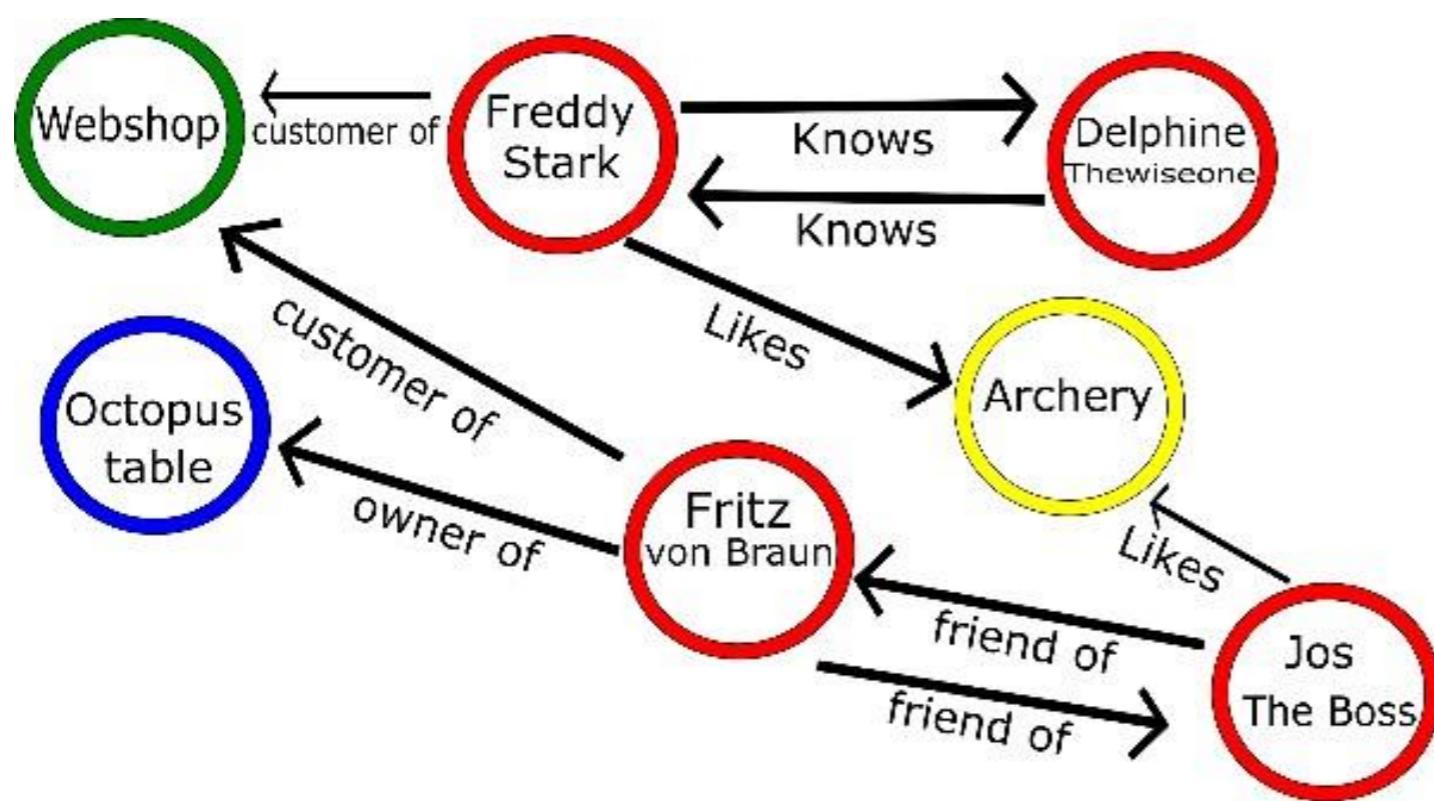
```
{
  "articles": [
    {
      "title": "title of the article",
      "articleID": 1,
      "body": "body of the artricle",
      "author": "Isaac Asimov",
      "comments": [
        {
          "username": "Fritz",
          "join date": "1/4/2014",
          "commentid": 1,
          "body": "this 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"
        }
      ]
    }
  ]
}
```

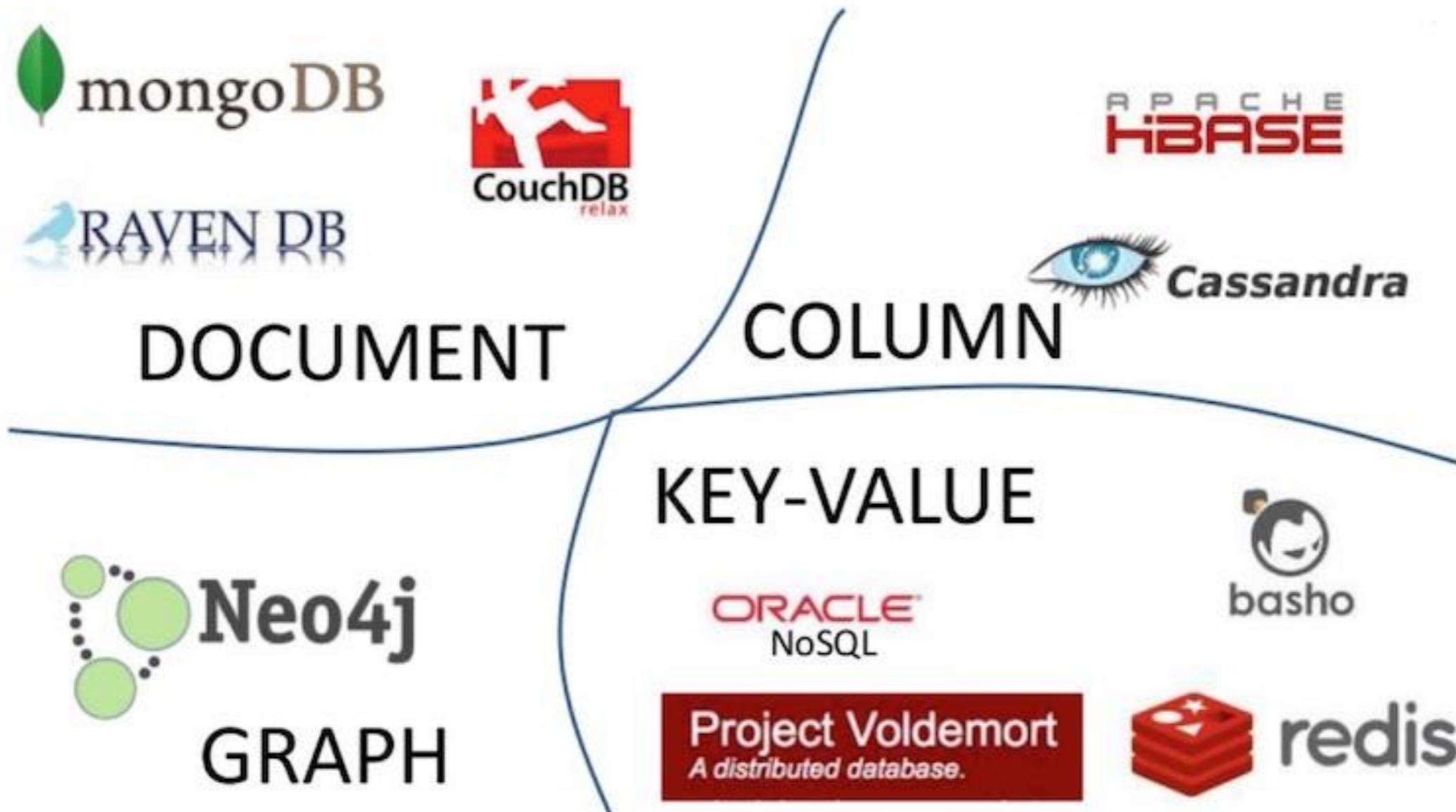
Whereas relational databases chop up data, Document stores save documents as a single entity

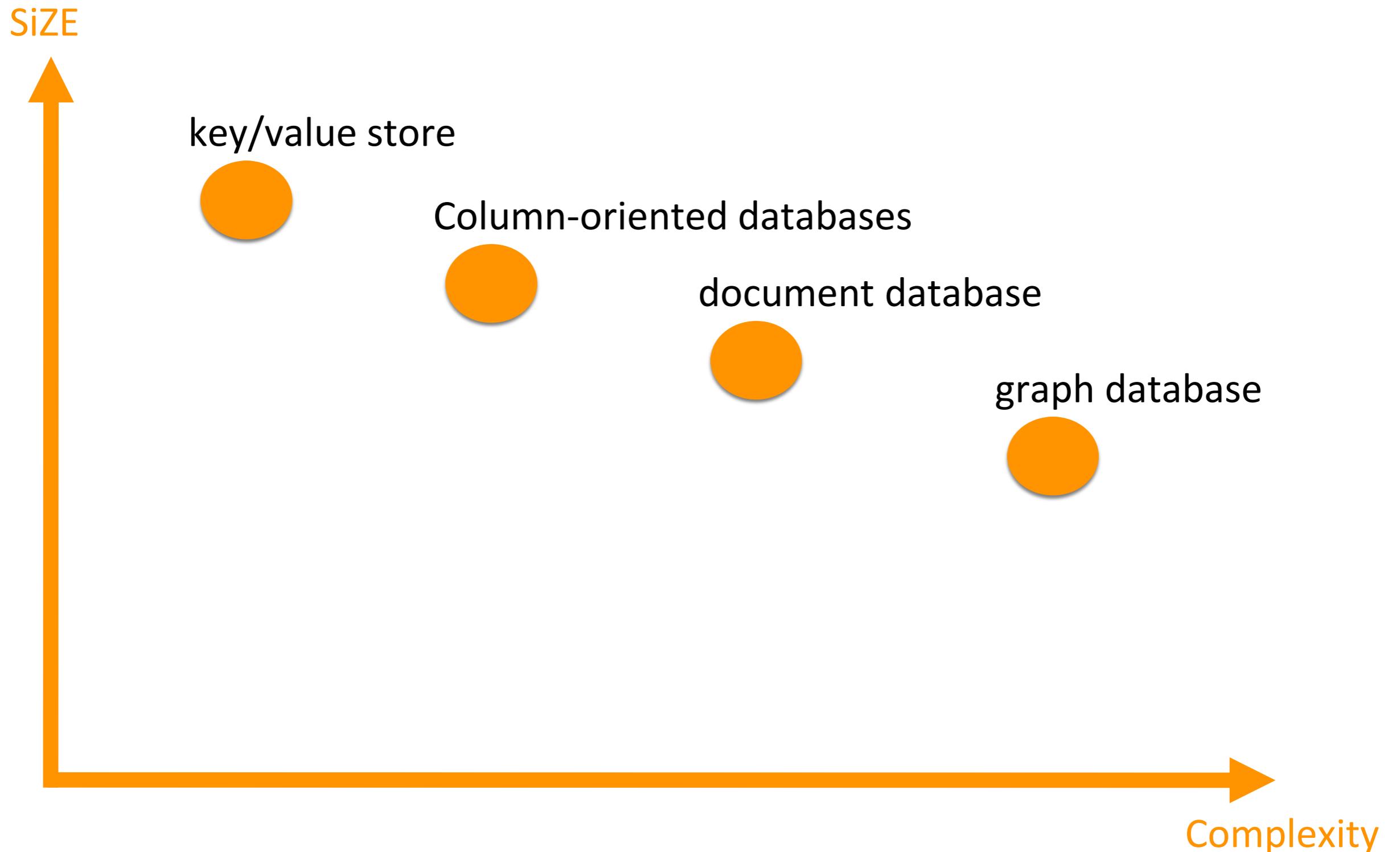


Graph database

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







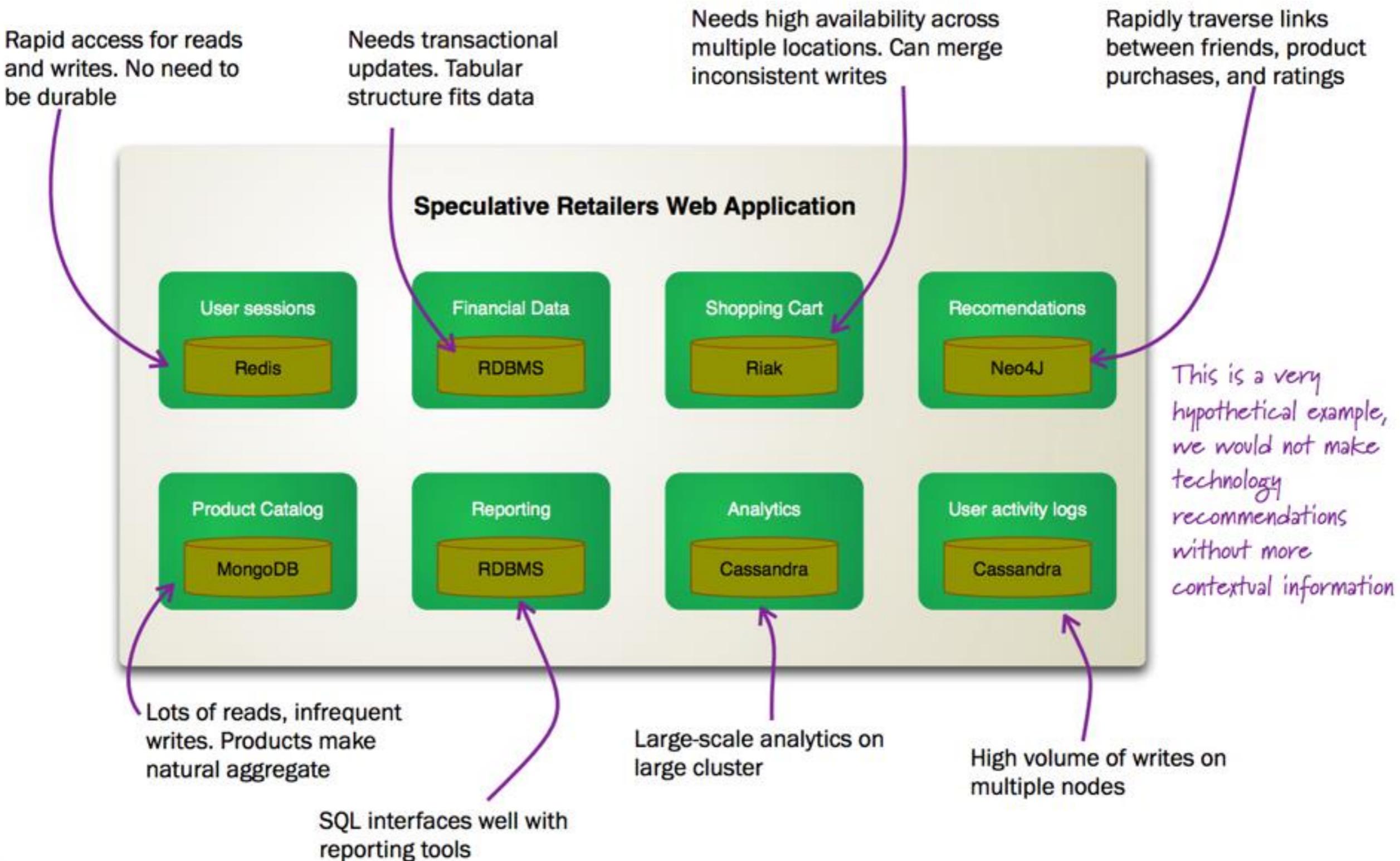
SQL

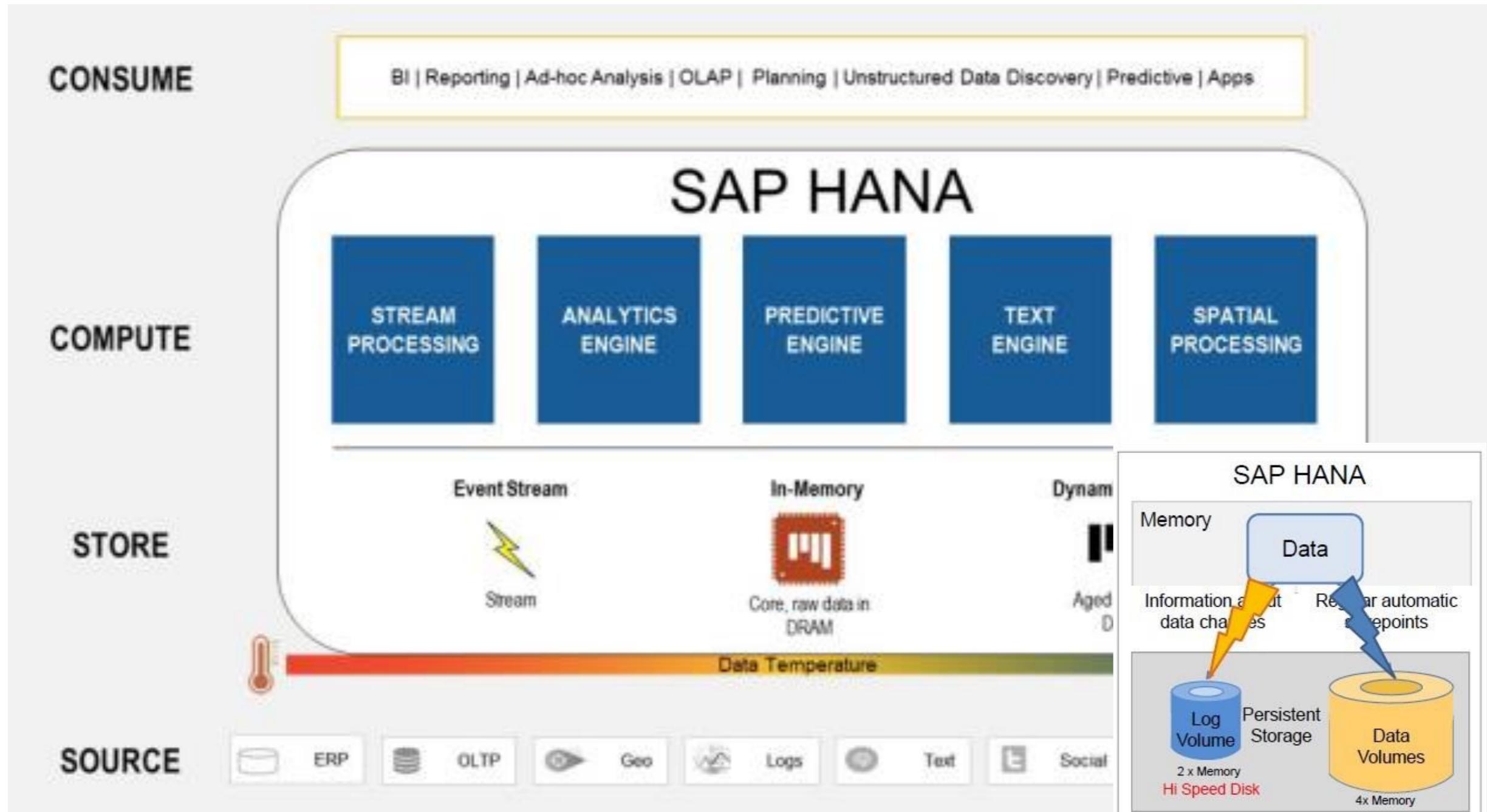
works great, can't scale for large data

NoSQL

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 \(PDF\)](#) ➤

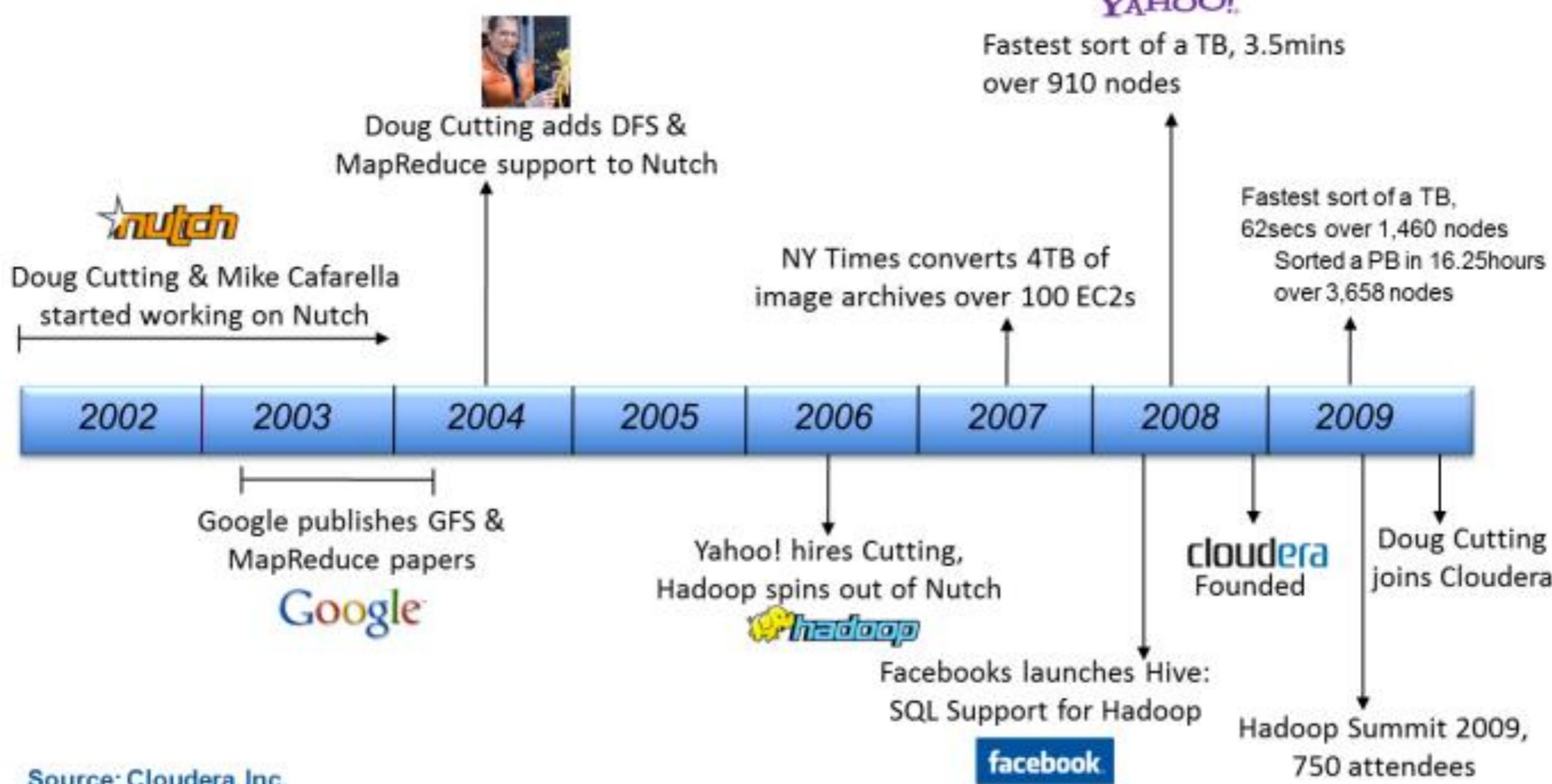
What is Hadoop ?

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

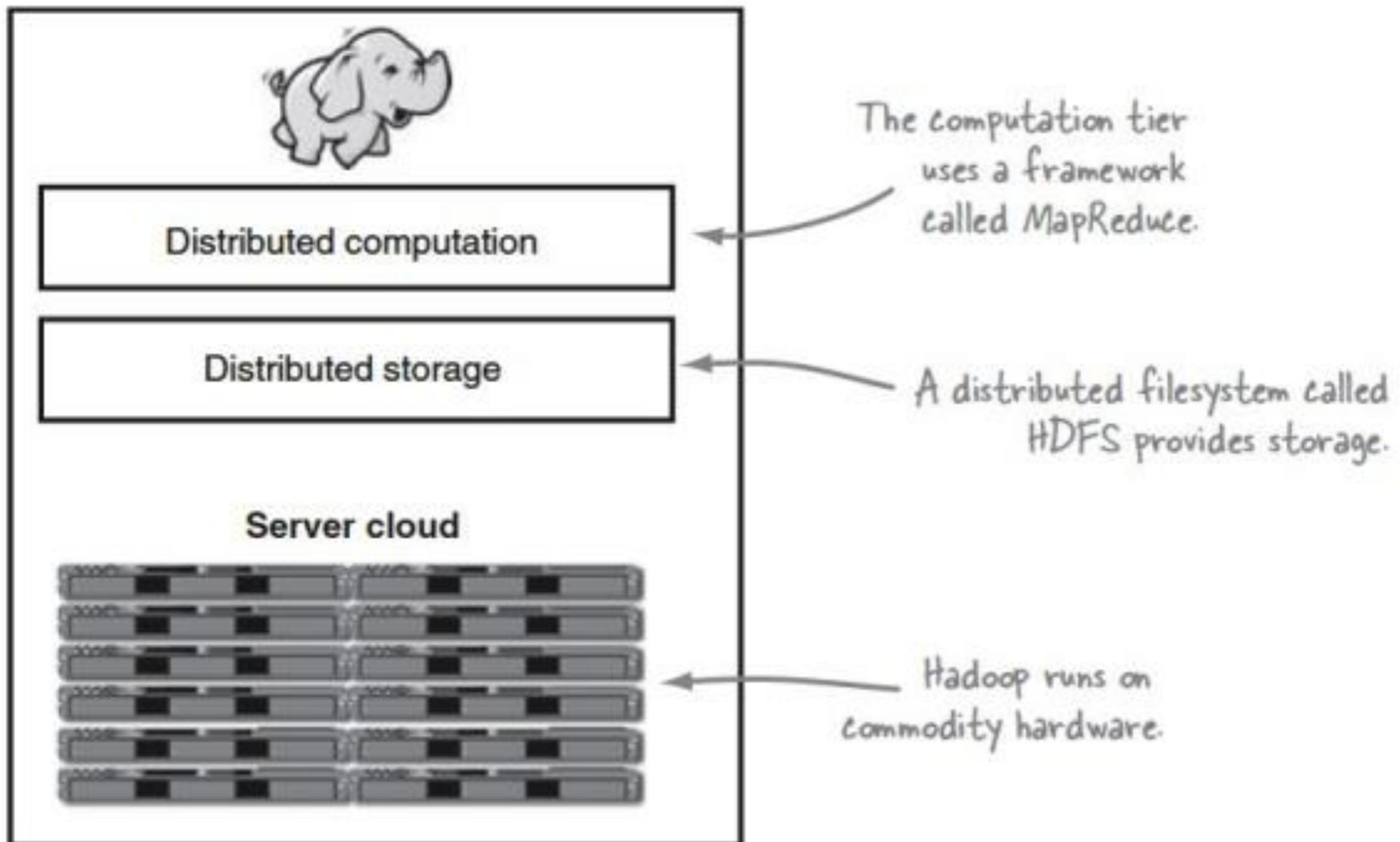
**Completely written in java
Open source & distributed under Apache license**



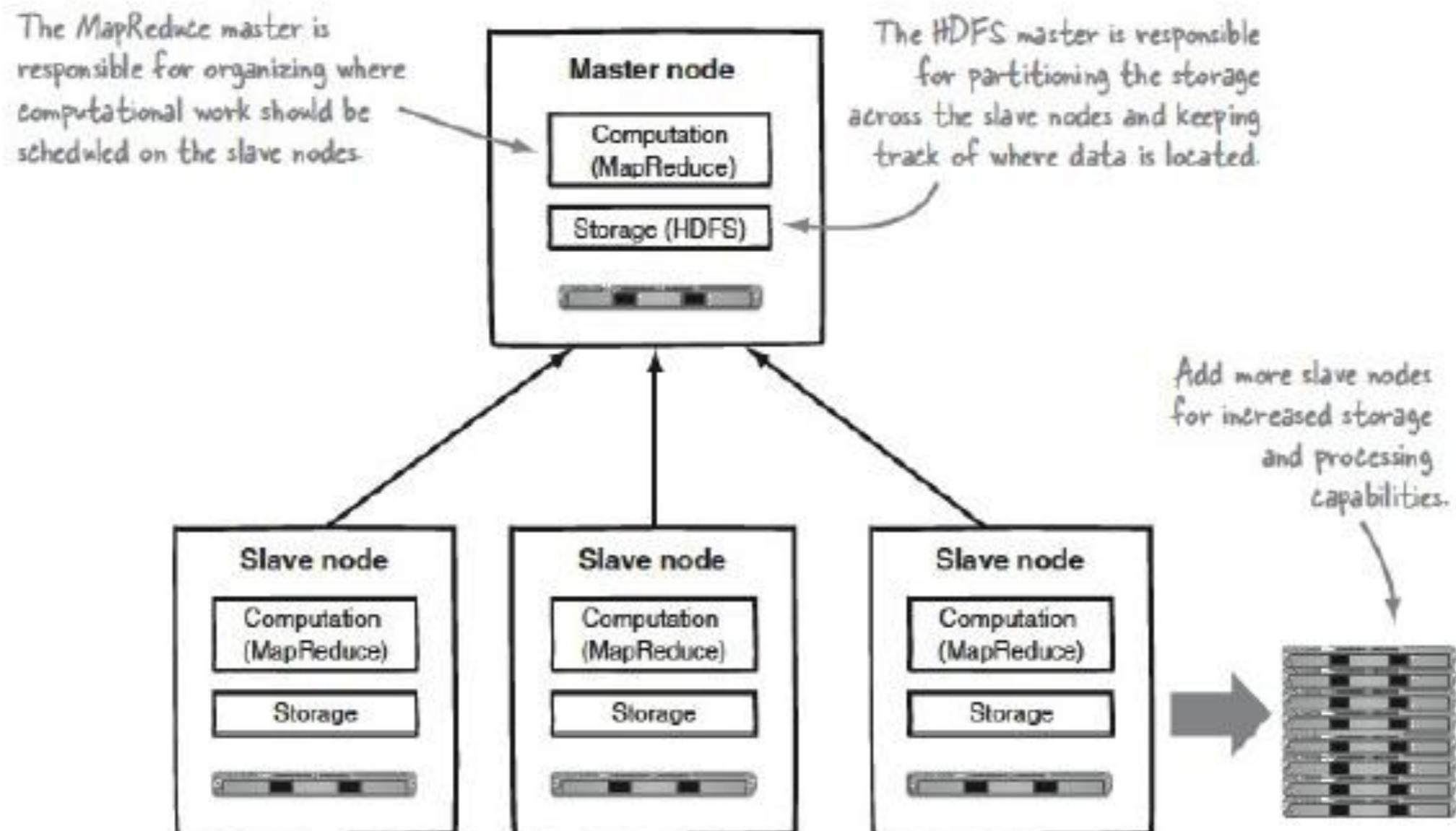
Hadoop Creation History



Hadoop Environment



Hadoop Architecture



Major Hadoop Components



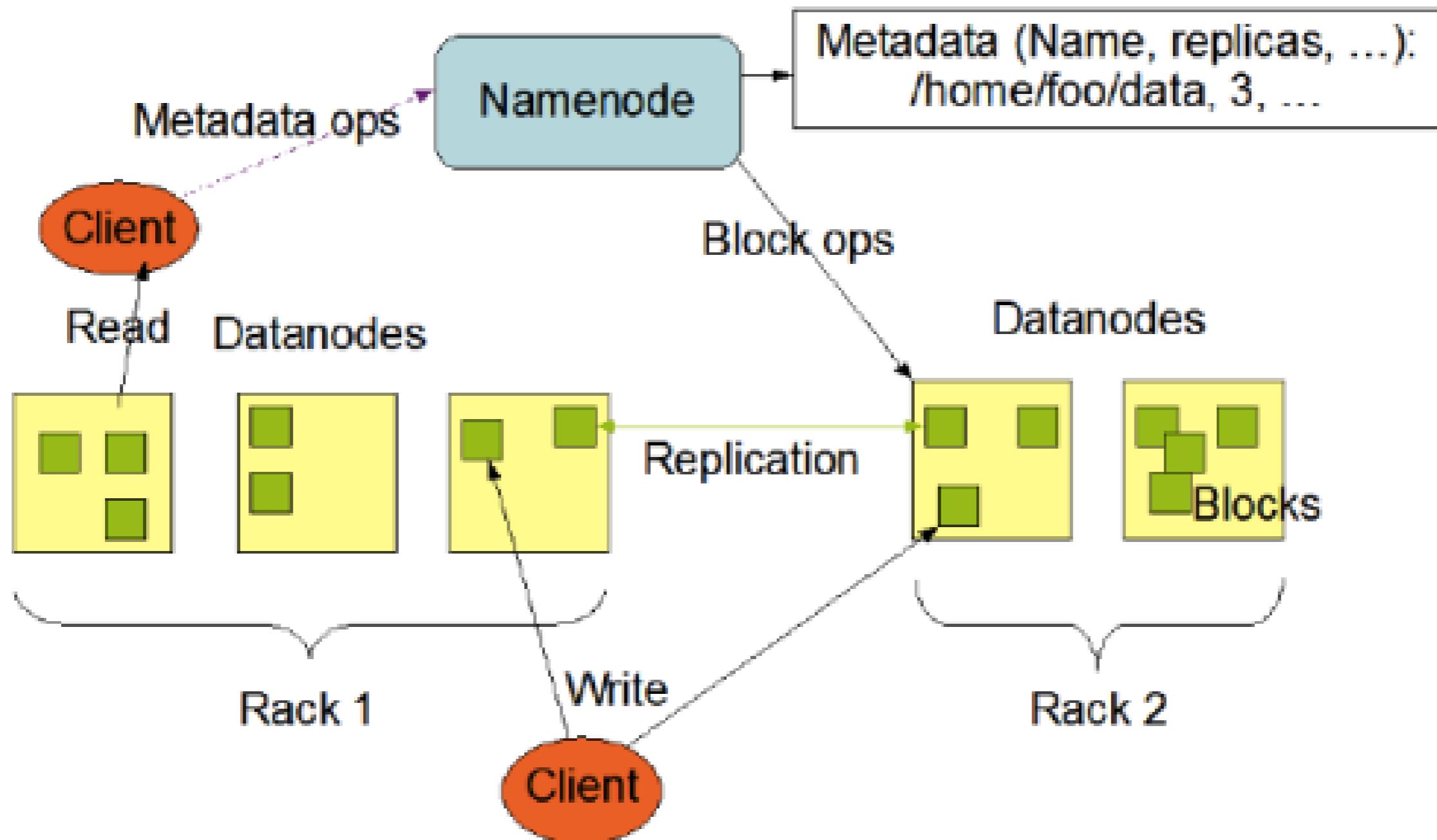
Hadoop Distributed File System (HDFS)

Map/Reduce System



- Default storage for the Hadoop cluster
- Data is distributed and replicated over multiple machines
- Designed to handle very large files with streaming data access patterns.
- NameNode/DataNode
- Master/slave architecture (1 master 'n' slaves)
- Designed for large files (64 MB default, but configurable) across all the nodes

HDFS Architecture

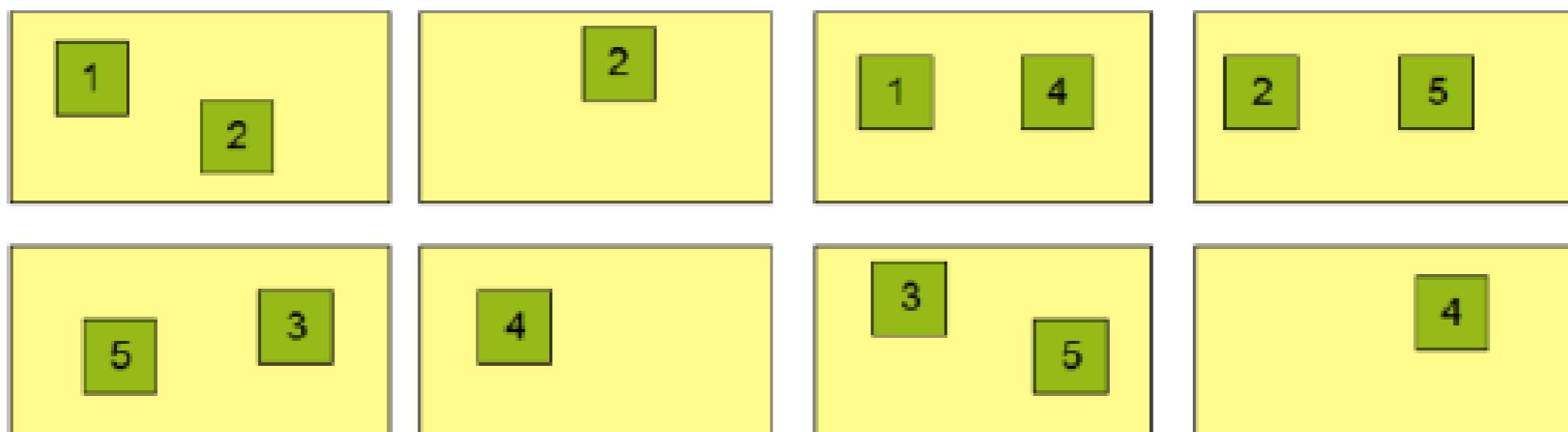


Data Replication in HDFS

Block Replication

```
Namenode (Filename, numReplicas, block-ids, ...)  
/users/sameerp/data/part-0, r:2, {1,3}, ...  
/users/sameerp/data/part-1, r:3, {2,4,5}, ...
```

Datanodes



How does HDFS work?

A file we want to store on HDFS ...

600 MB

We're raising the question because no one else wants to, because no one else wants to say what needs to be said.

And let's be real, it's the two-ton elephant in the room with nearly every other star's name on the trade rumor radar these days.

We've read over and over again about Nash refusing to ask for a trade, refusing to play the game that so many others have late in their careers.

How does HDFS work?

HDFS Splits file into blocks ...

256 MB

We're raising the question because no one else wants to, because no one else wants to say what needs to be said.

256 MB

And let's be real, it's the two-ton elephant in the room with nearly every other star's name on the trade rumor radar these days.

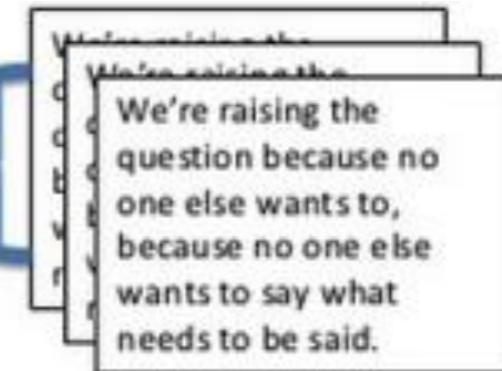
88 MB

We've read over and over again about Nash refusing to play the game that so many others have late in their careers.

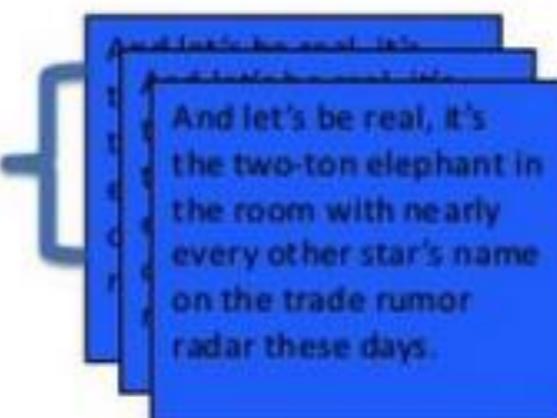
How does HDFS work?

HDFS will create **3replicas** of each block ...

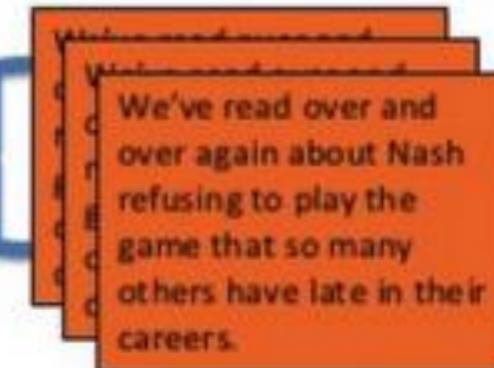
3 copies



3 copies

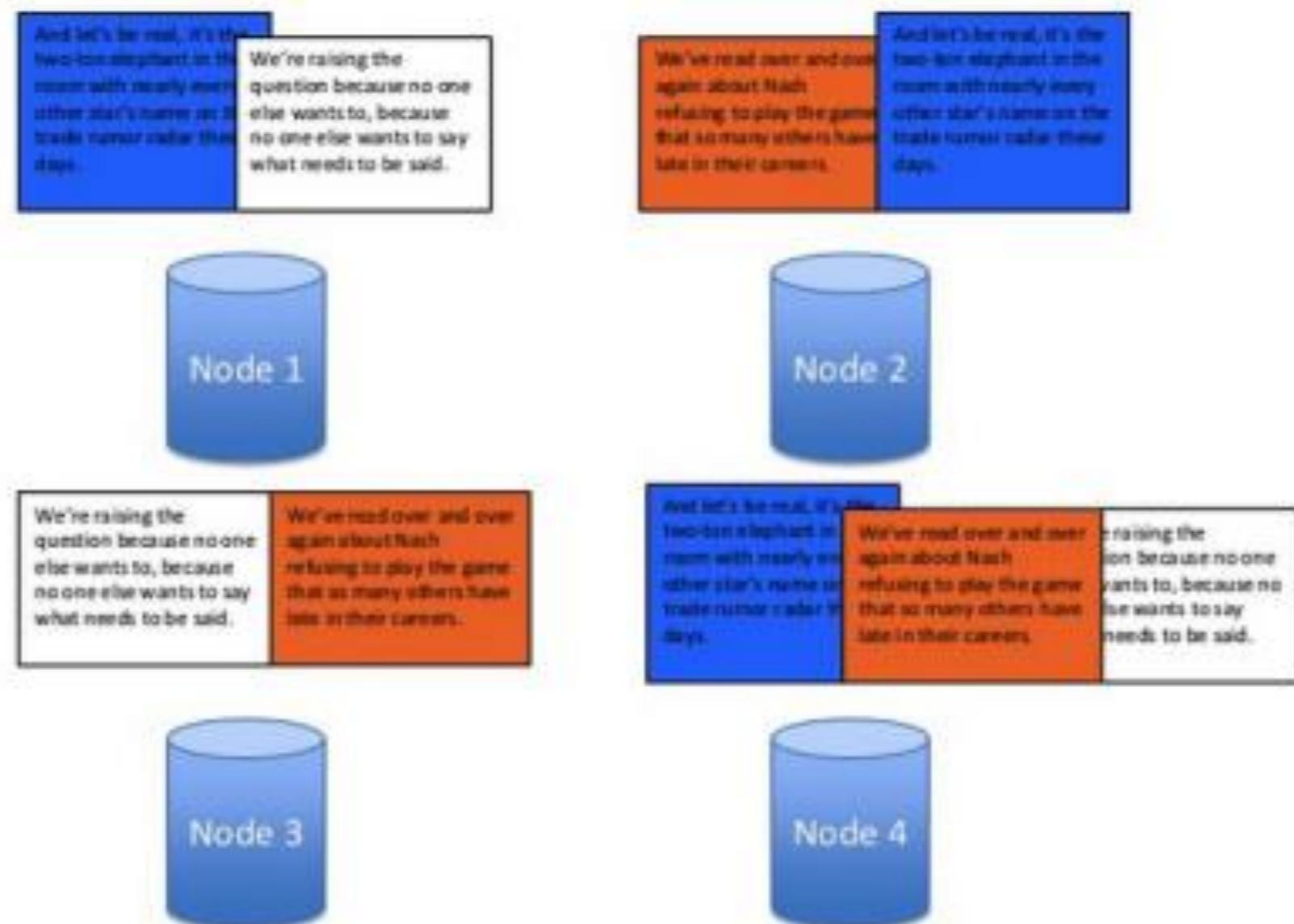


3 copies



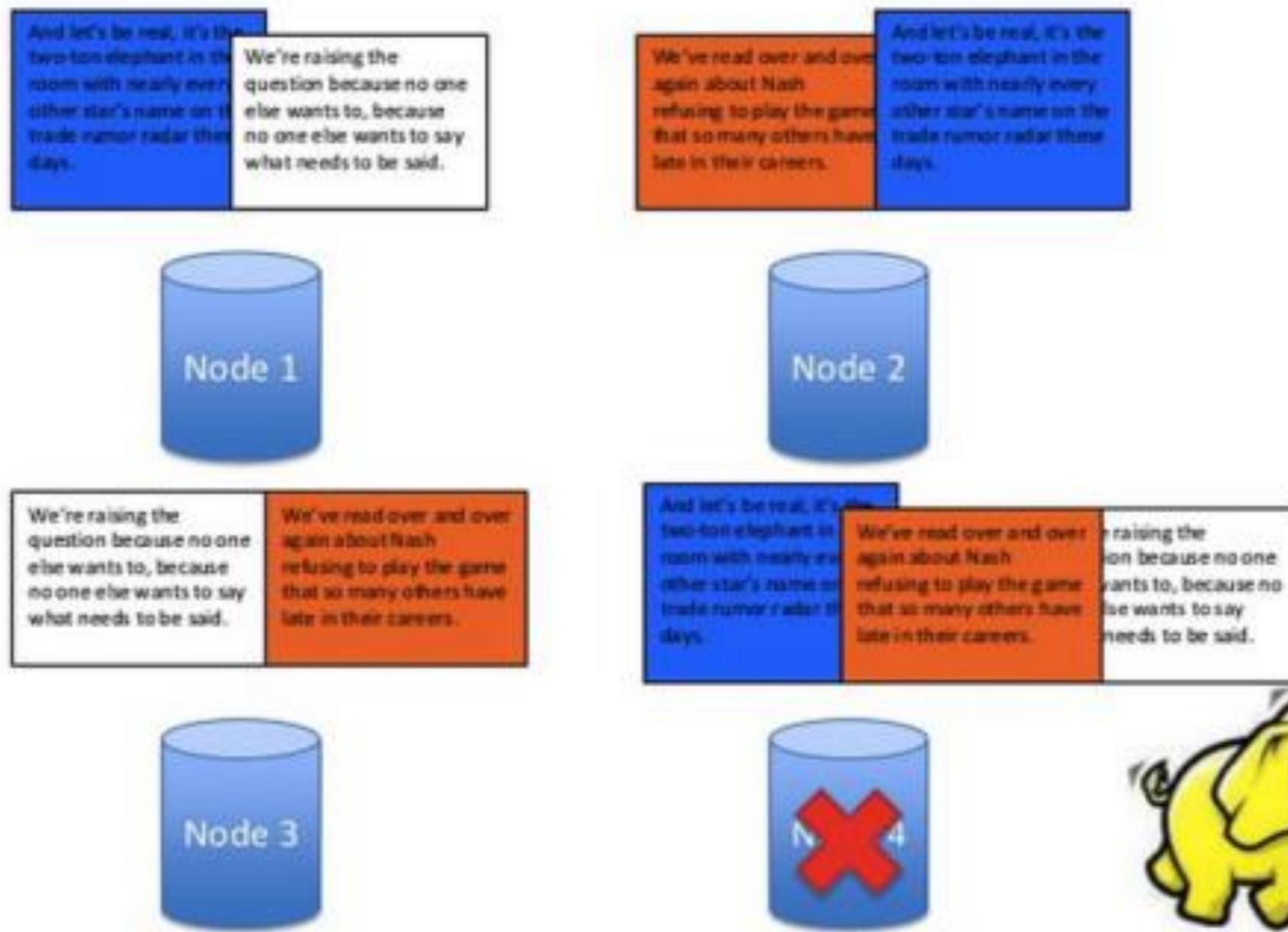
How does HDFS work?

HDFS distributes these replicas across the cluster ...

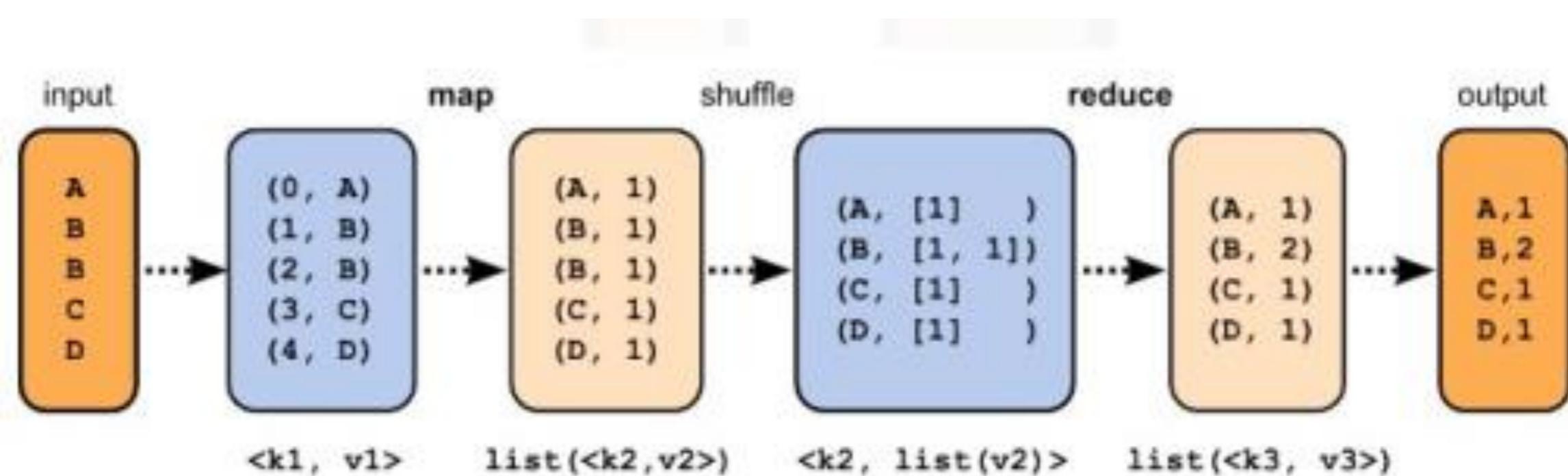


How does HDFS work?

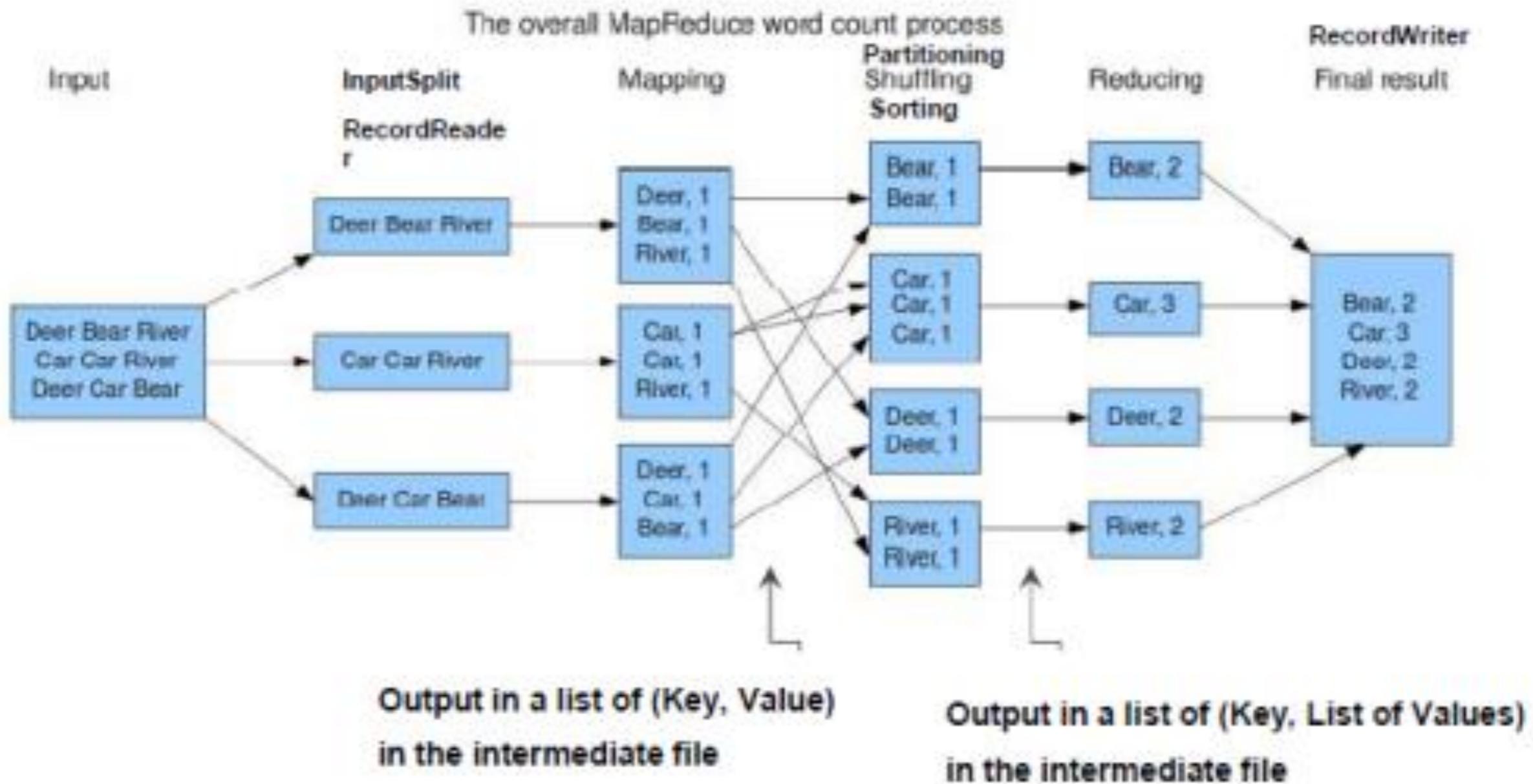
If a node goes down, we have copies elsewhere



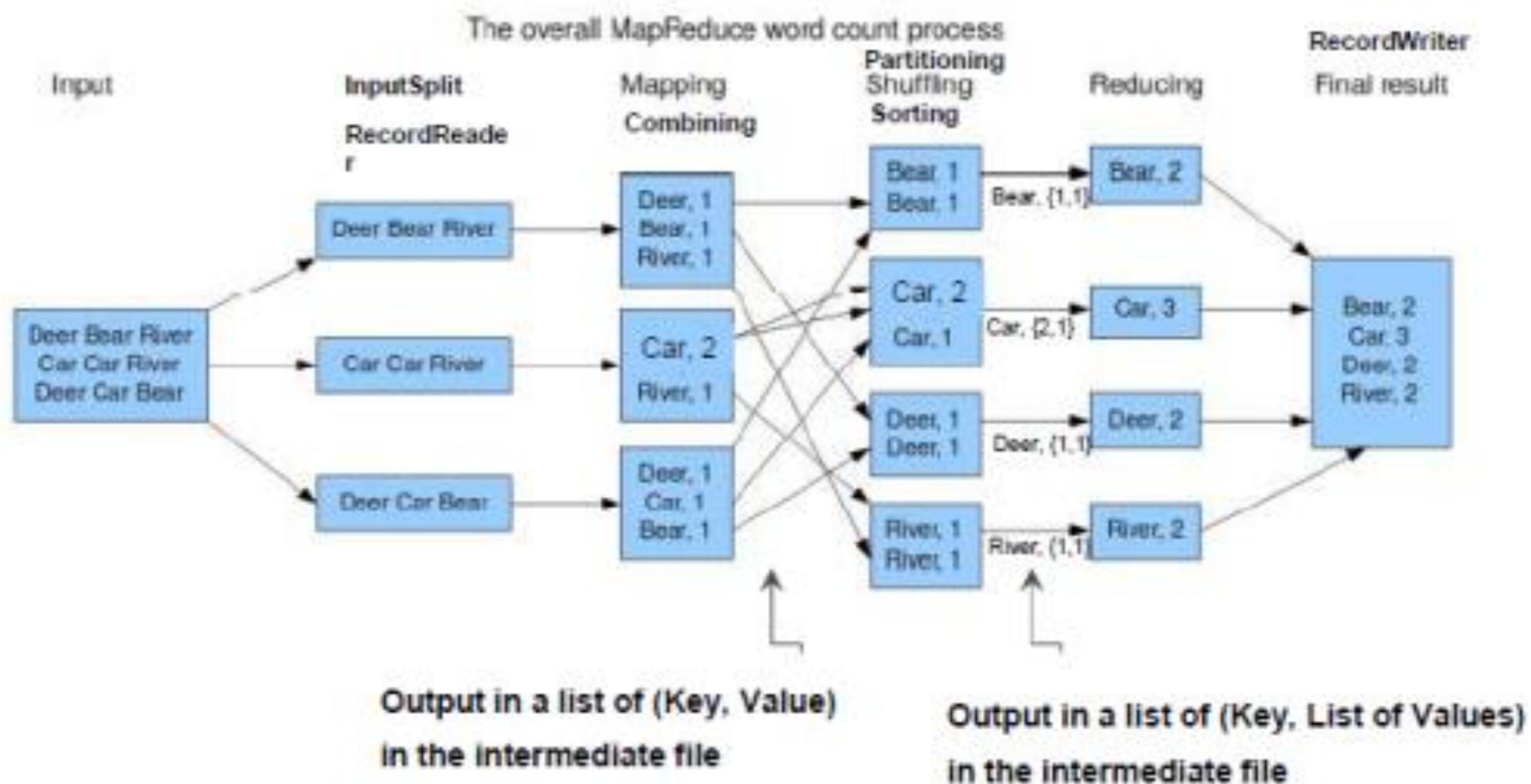
MapReduce Framework



How does the MapReduce work ?



How does the MapReduce work ?



WordCount - MapReduce

```
public class WordCount {  
  
    public static class Map extends MapReduceBase implements  
        Mapper<LongWritable, Text, Text, IntWritable> {  
        private final static IntWritable one = new IntWritable(1);  
        private Text word = new Text();  
  
        public void map(LongWritable key, Text value, OutputCollector<Text, IntWritable>  
            ] output, Reporter reporter) throws IOException {  
            String line = value.toString();  
            StringTokenizer tokenizer = new StringTokenizer(line);  
            while (tokenizer.hasMoreTokens()) {  
                word.set(tokenizer.nextToken());  
                output.collect(word, one);  
            }  
        }  
  
        public static class Reduce extends MapReduceBase implements  
            Reducer<Text, IntWritable, Text, IntWritable> {  
            public void reduce(Text key, Iterator<IntWritable> values, OutputCollector<Text,  
                IntWritable> output, Reporter reporter) throws IOException {  
                int sum = 0;  
                while (values.hasNext()) { sum += values.next().get(); }  
                output.collect(key, new IntWritable(sum));  
            }  
        }  
    }  
}
```

Map Function

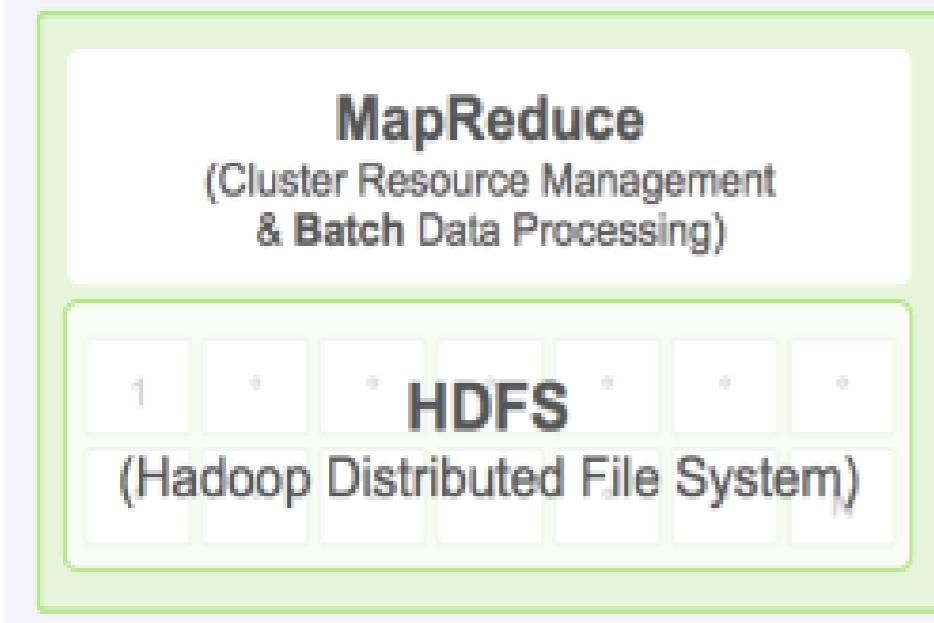
Reduce Function

Hadoop 2.X



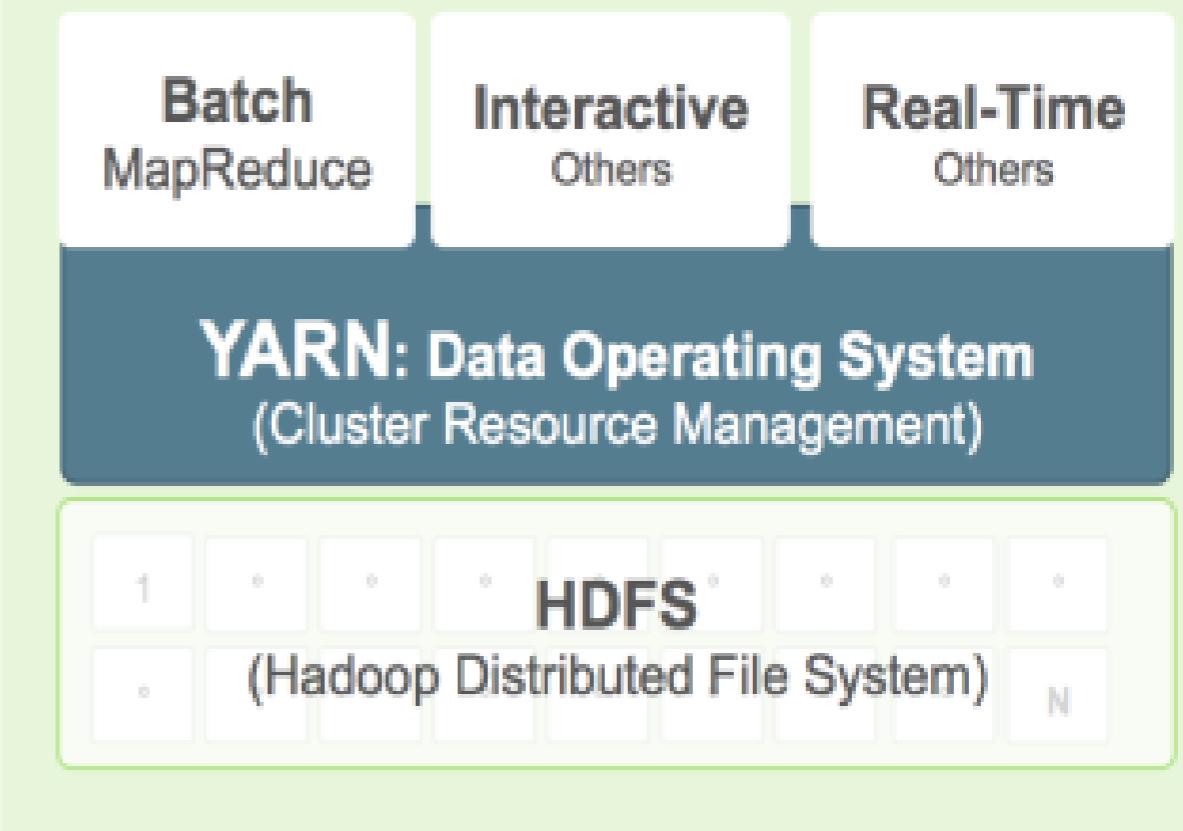
Hadoop 1

- Silos & Largely batch
 - Single Processing engine

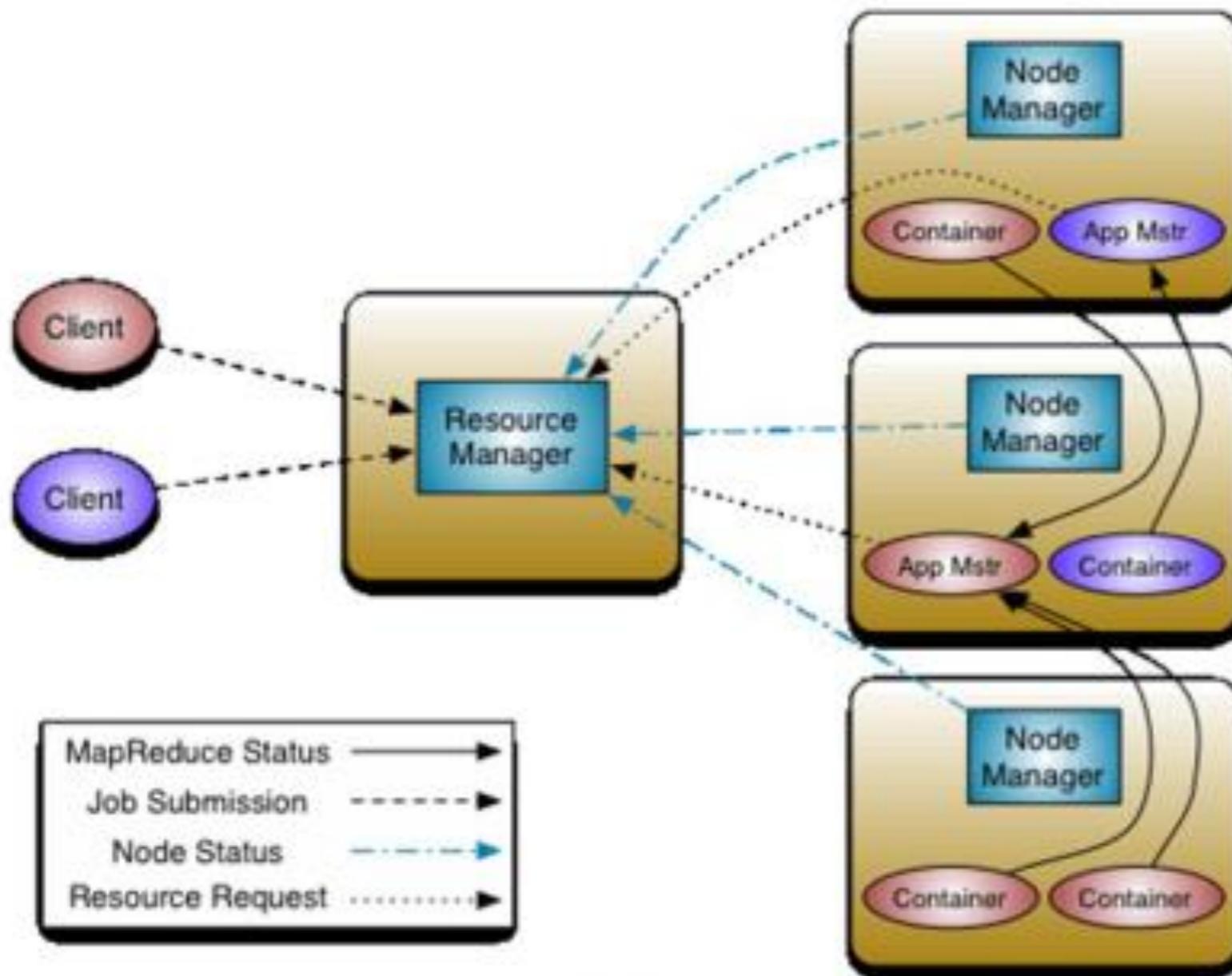


Hadoop 2 w/YARN

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



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

Evolution of the Hadoop Platform

The stack is continually evolving and growing!

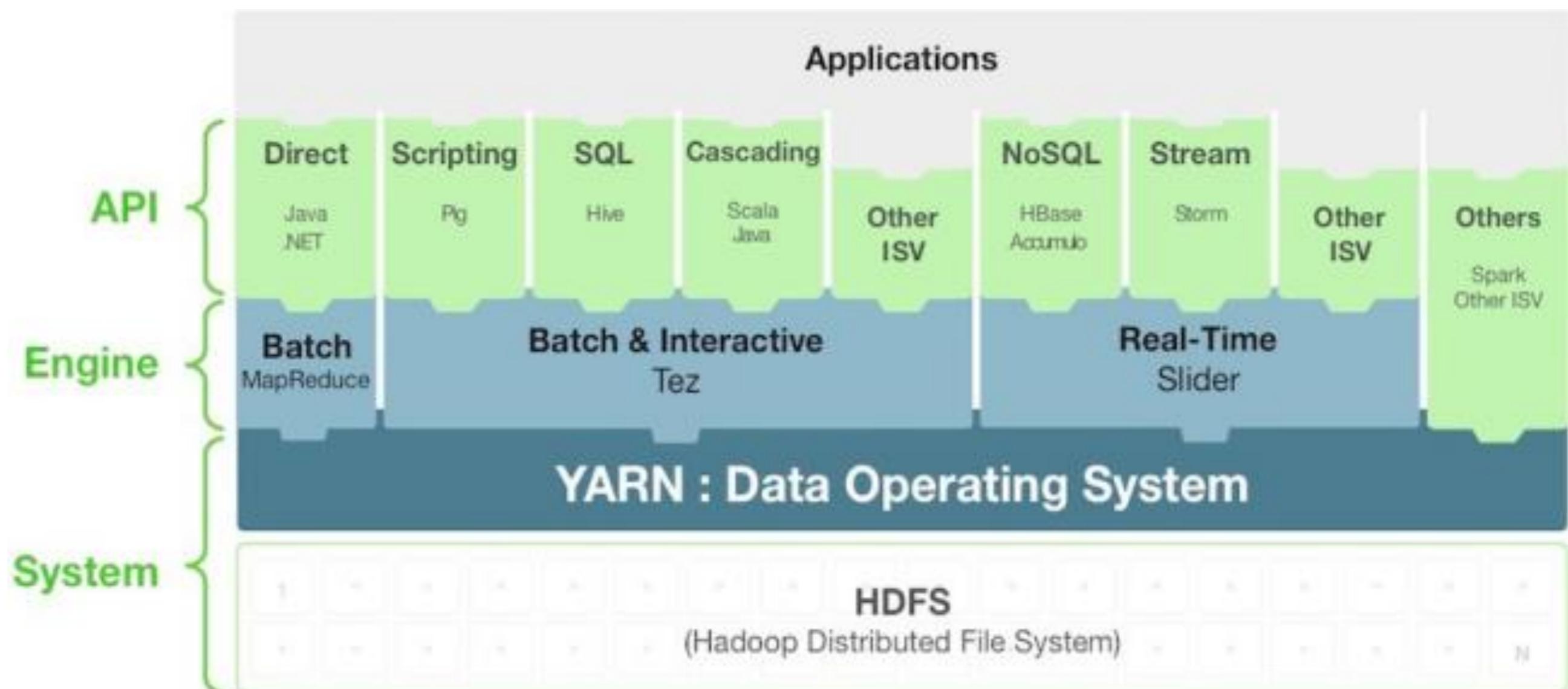
Core Hadoop (HDFS, MapReduce)	Solr Pig	Core Hadoop	HBase ZooKeeper	Solr Pig	Core Hadoop	Hive Mahout	Hive Mahout	Hive Mahout	Hive Mahout	Ibis Flink Parquet
2006	2007	2008	2009	2010	2011	2012	2013	2014-15		
Core Hadoop (HDFS, MapReduce)	Solr Pig	Core Hadoop	HBase ZooKeeper	Solr Pig	Core Hadoop	Hive Mahout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

cloudera

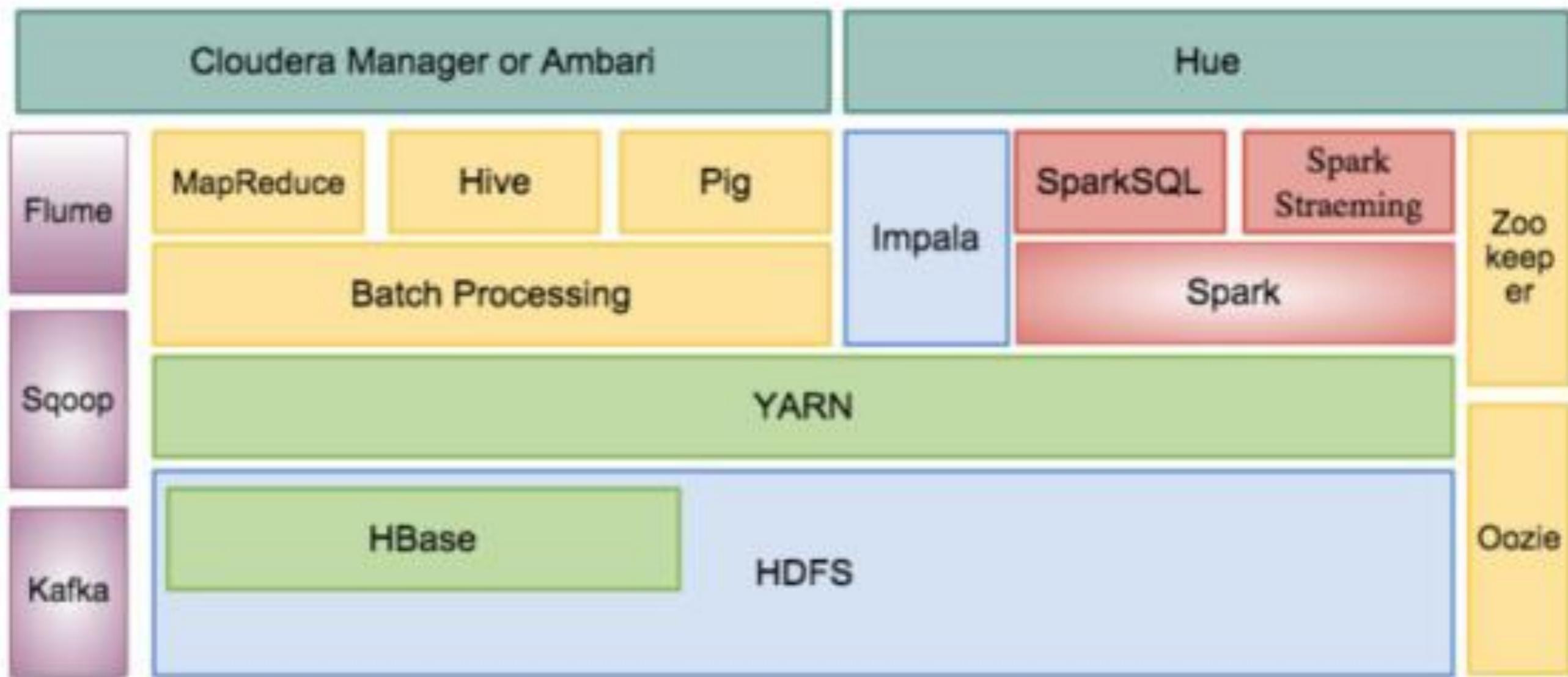
© Cloudera, Inc. All rights reserved.

9

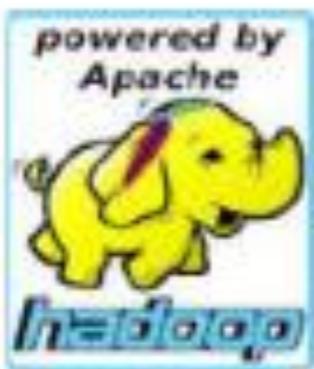
Hadoop 2.x Ecosystems



Hadoop Ecosystems



Hadoop Distribution



MAPR

cloudera



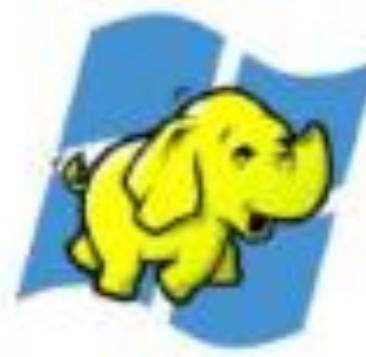
Pivotal



TERADATA

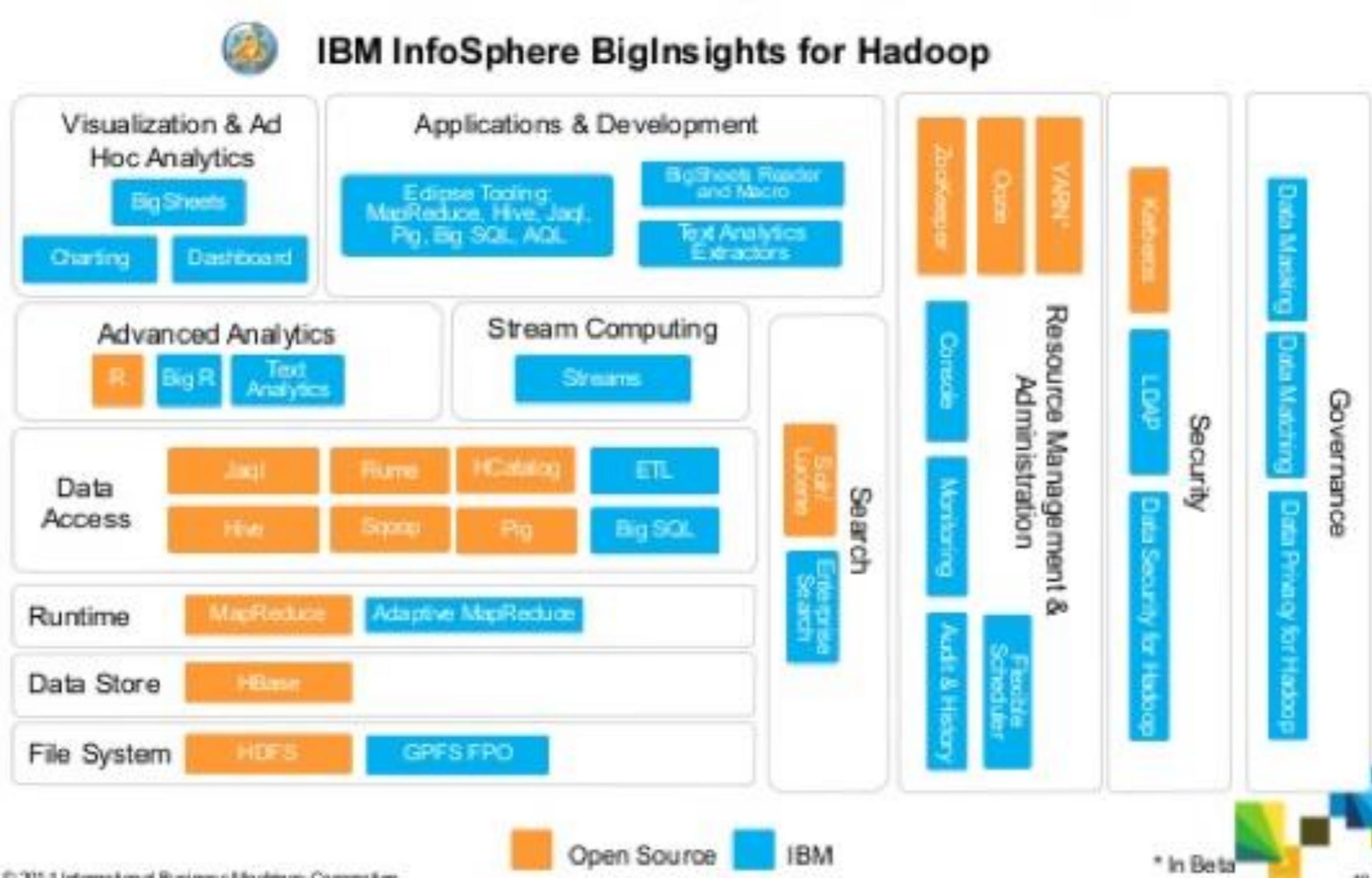


amazon
web services

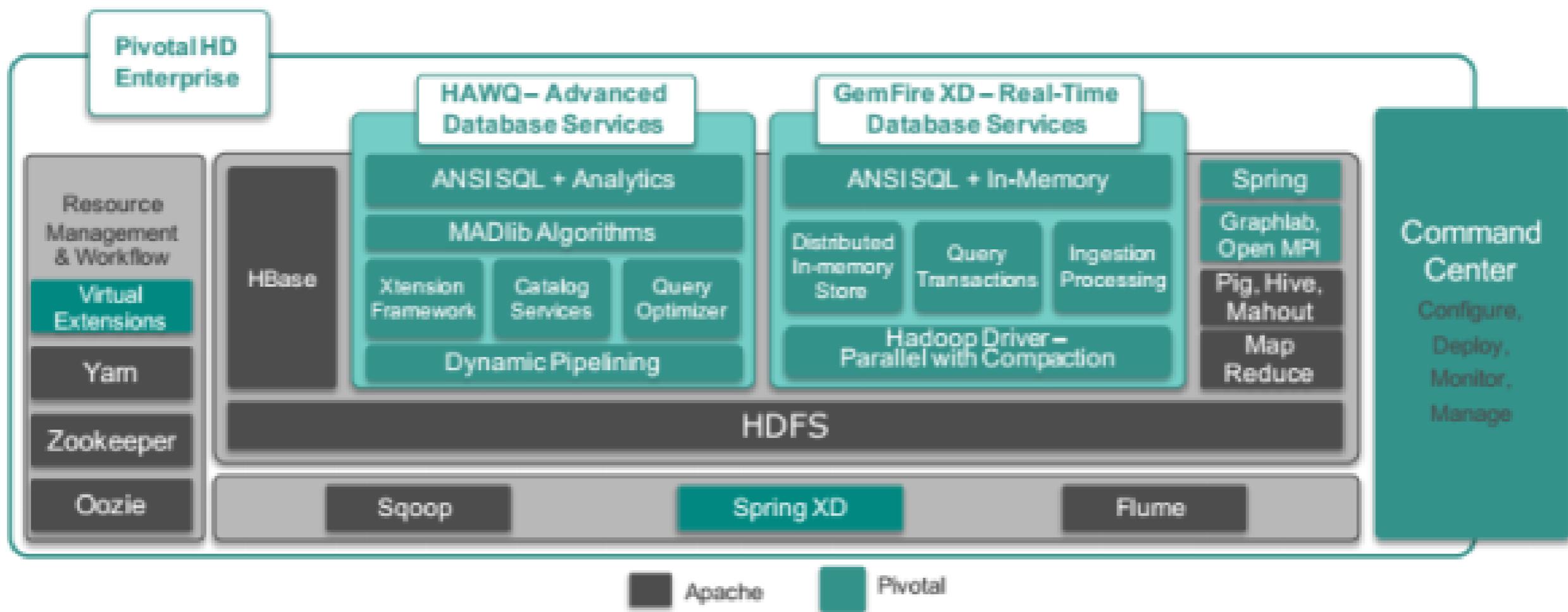


Microsoft Azure

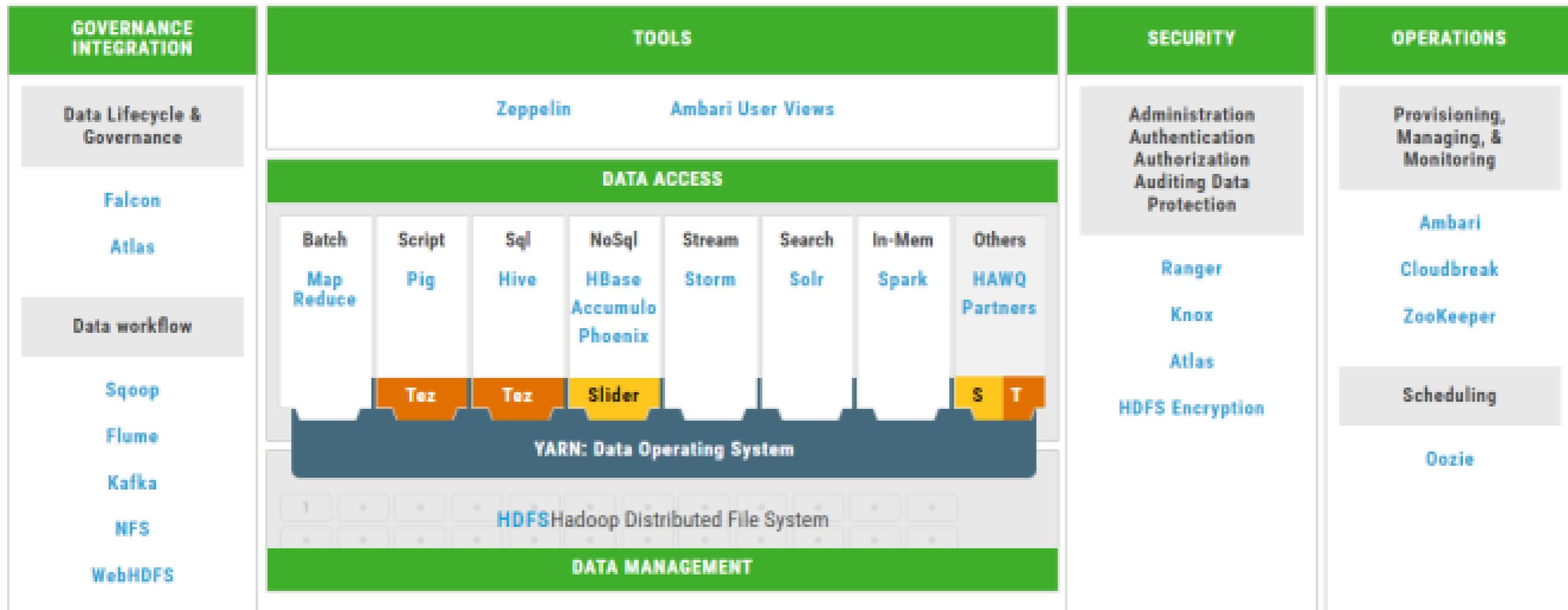
IBM InfoSphere BigInsights

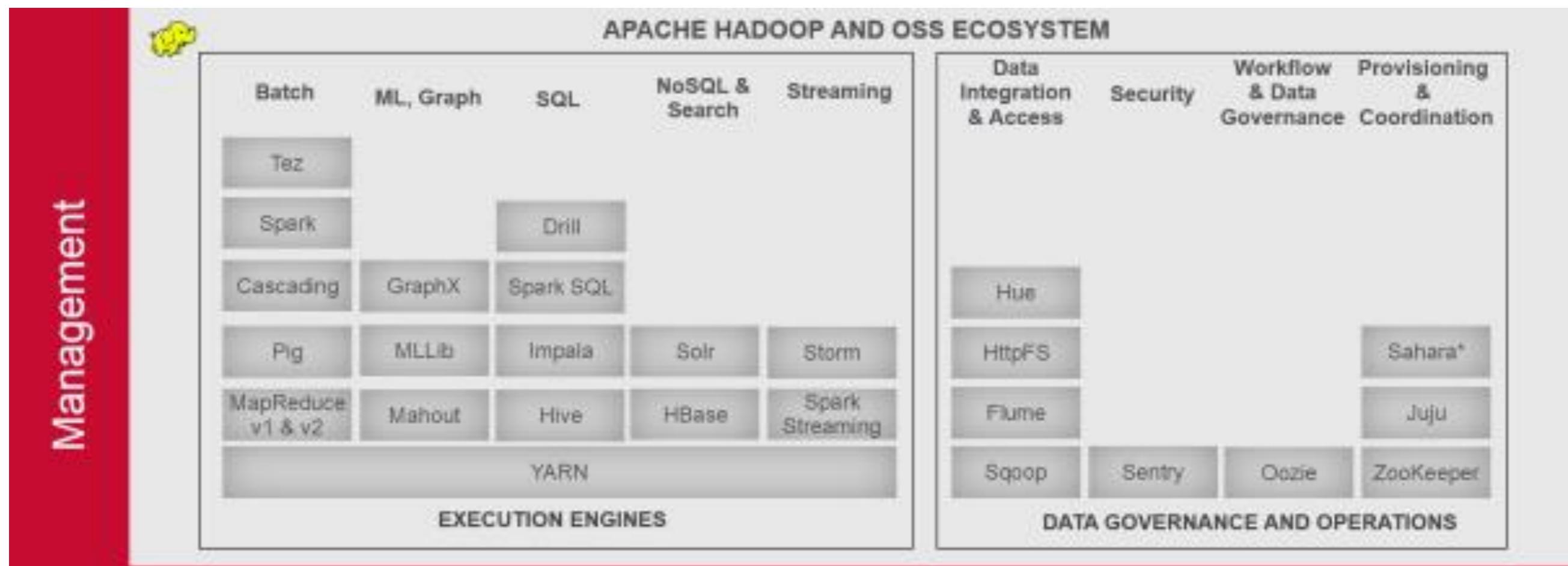


Pivotal HD Architecture



Hortonworks

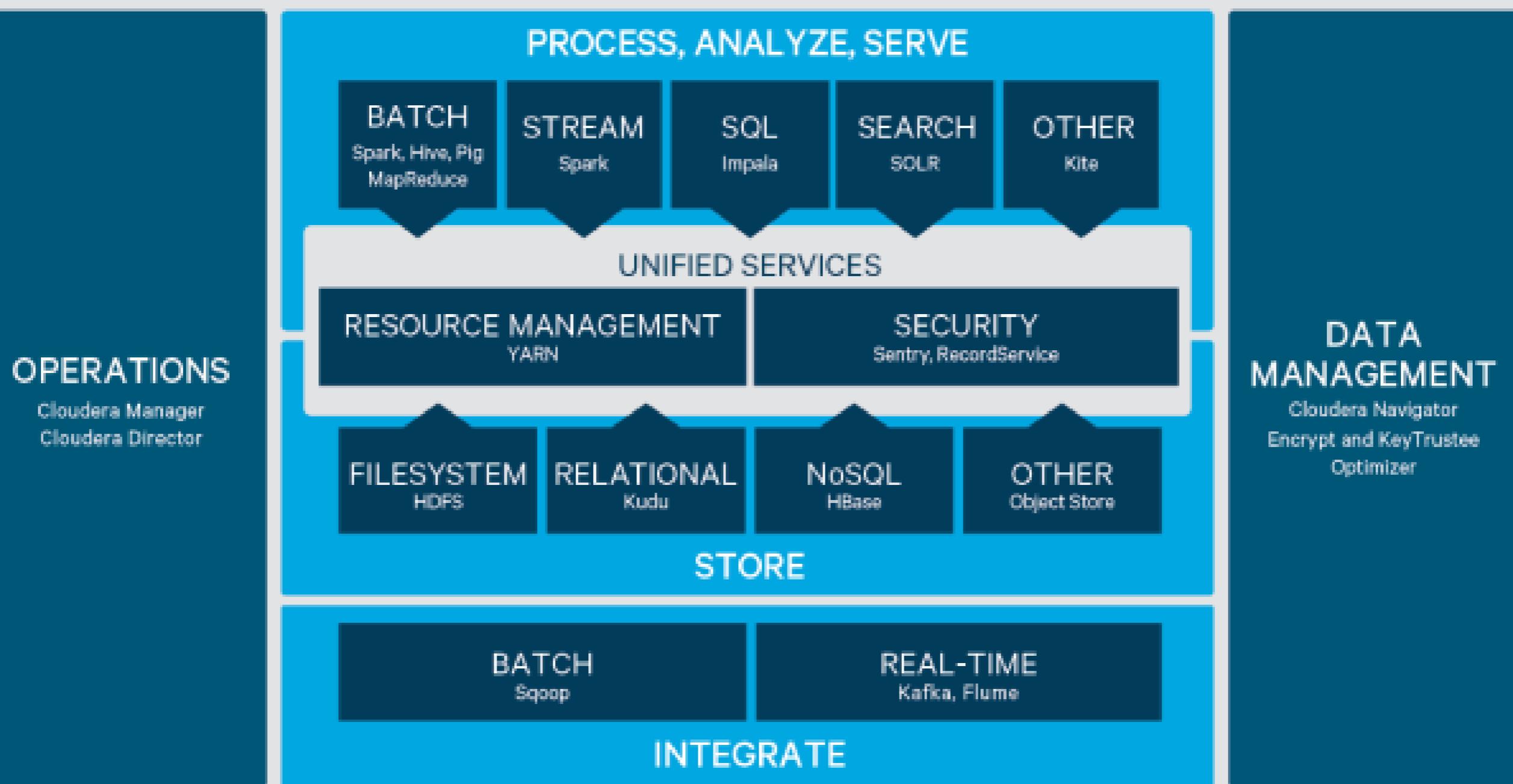




MapR-FS

Data Platform

MapR-DB



Default Cloudera Services

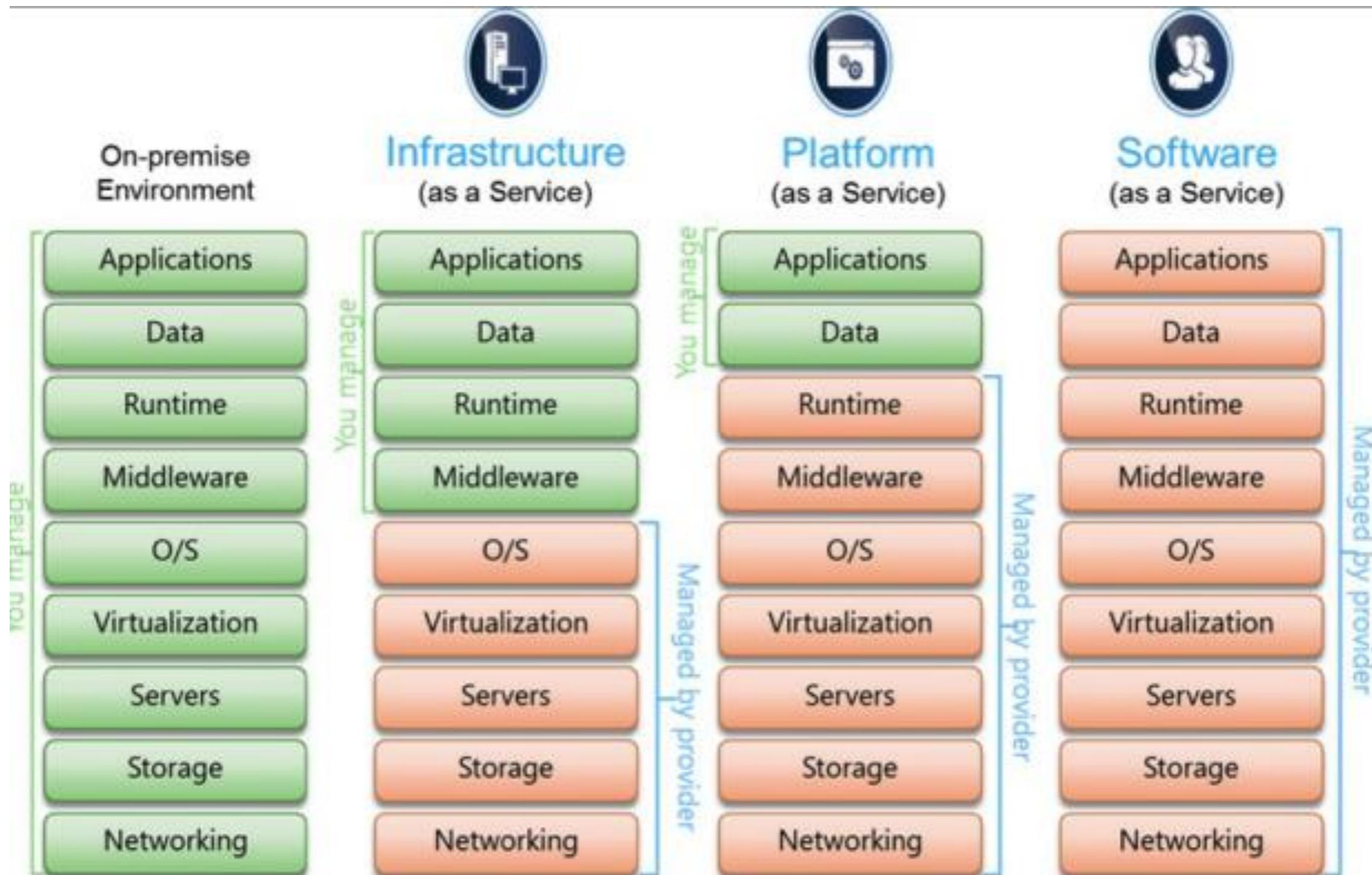
- Cloudera Manager
- HDFS
- YARN
- Apache Hive and Pig
- Apache Flume and Sqoop
- Apache Oozie
- Cloudera Hue
- ZooKeeper

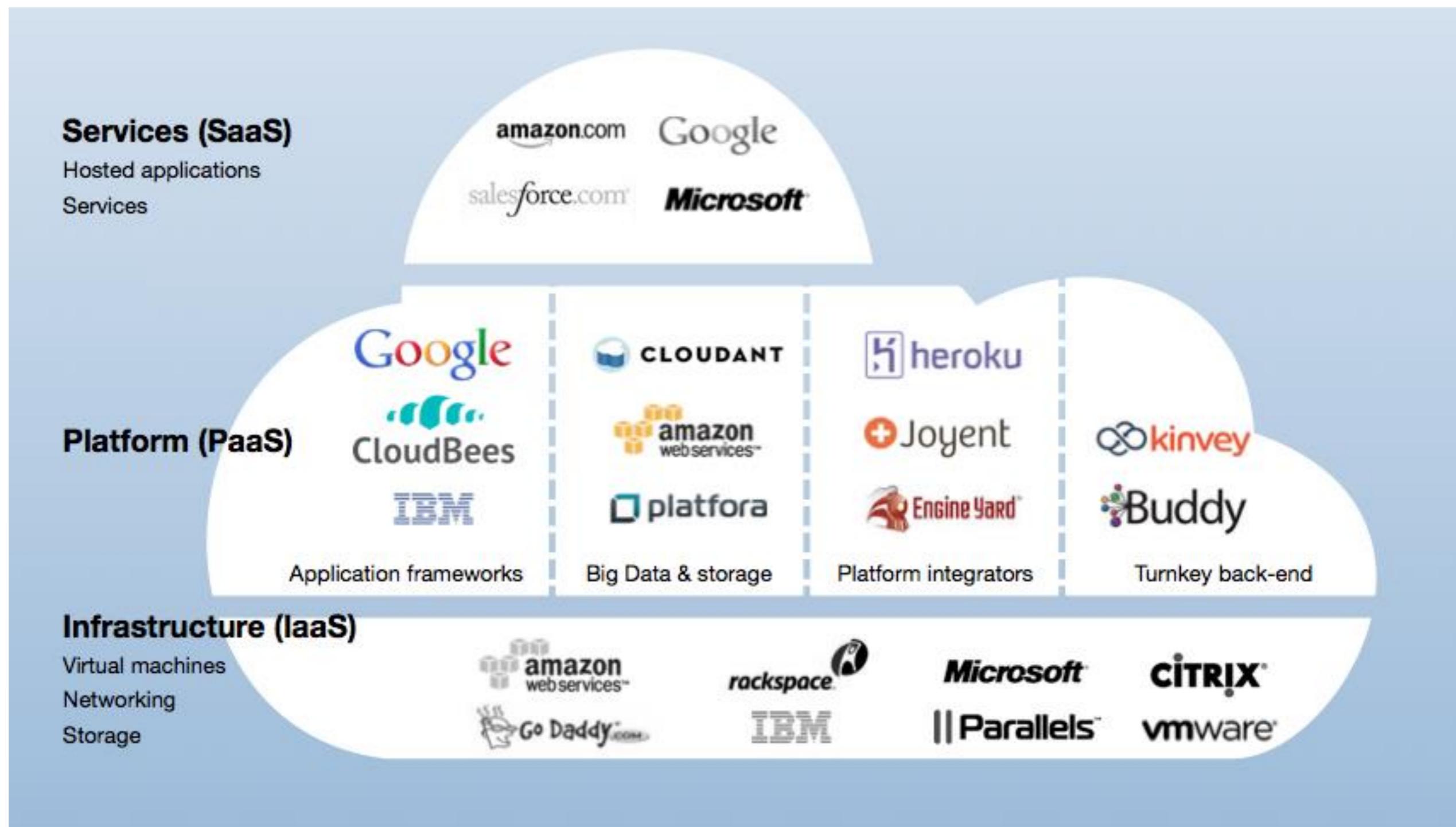


Issue with Big Data Infrastructure

- Large investment
- Scalability
- ROI
- Business Cases

Cloud Technology





Big Data on Cloud



Microsoft Azure

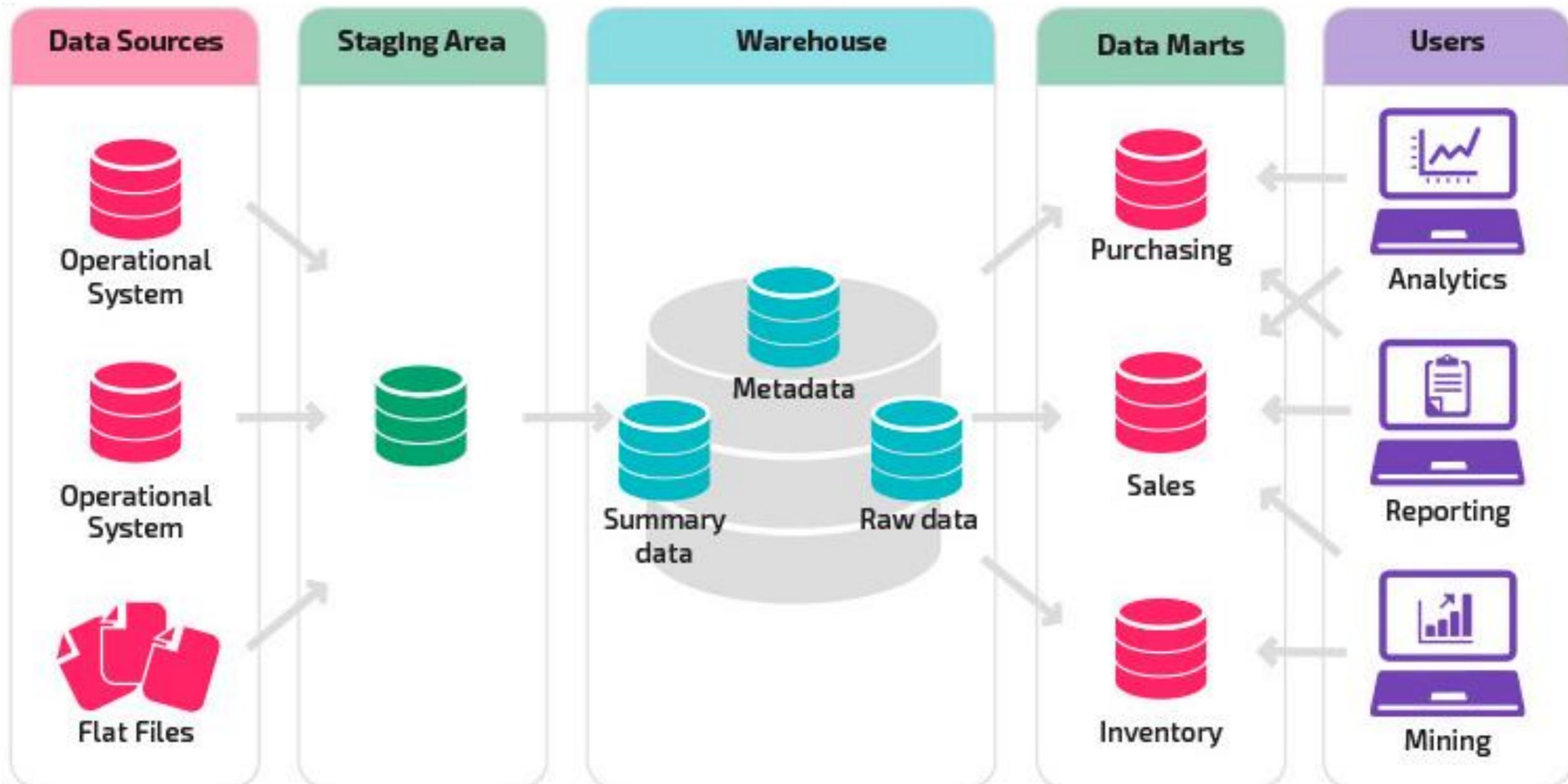


DATA LAKE

DATA LAKE



Data Warehouse



The old way: Ask, then collect

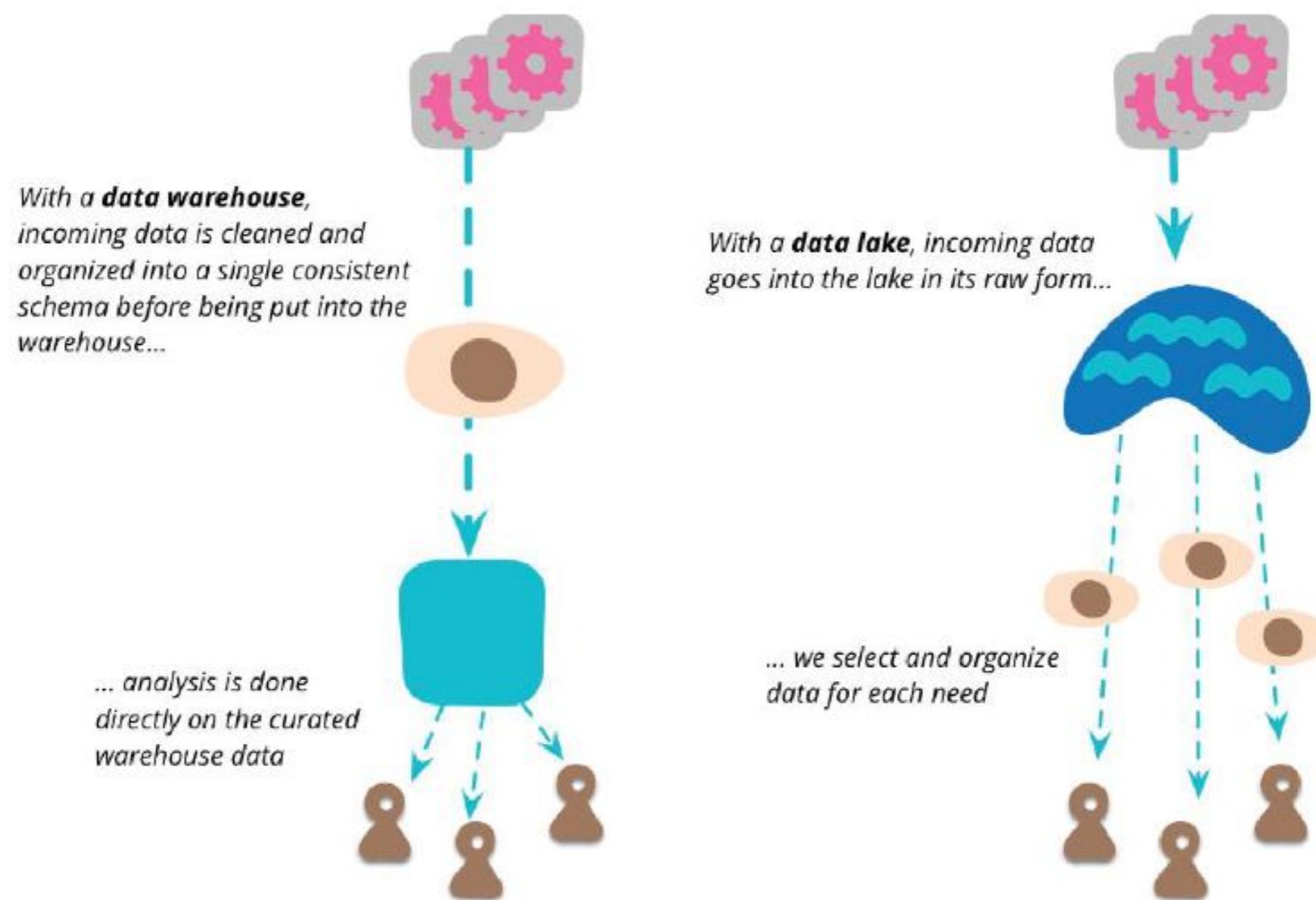


The new way: Collect, then ask

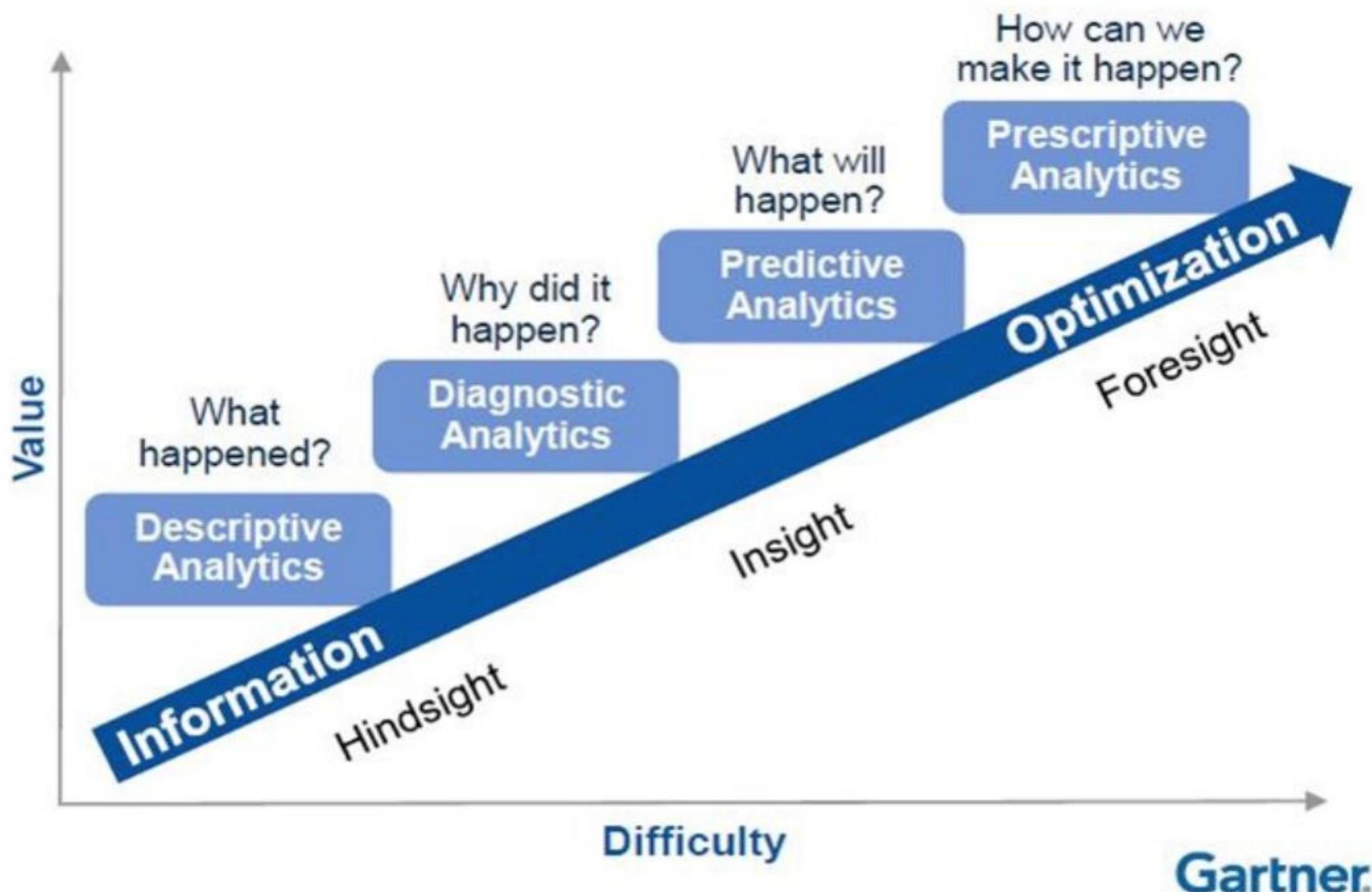


Source: Data Science and Critical Thinking, A. Croll

Differences between Data Lake and Data Warehouse



Source: martinfowler.com/bliki/DataLake.html



Traditional Analytics (BI)

vs

Big Data Analytics

Focus on

- Descriptive analytics
- Diagnosis analytics

Data Sets

- Limited data sets
- Cleansed data
- Simple models

Supports

Causation: what happened,
and why?

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

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

Correlation: new insight
More accurate answers

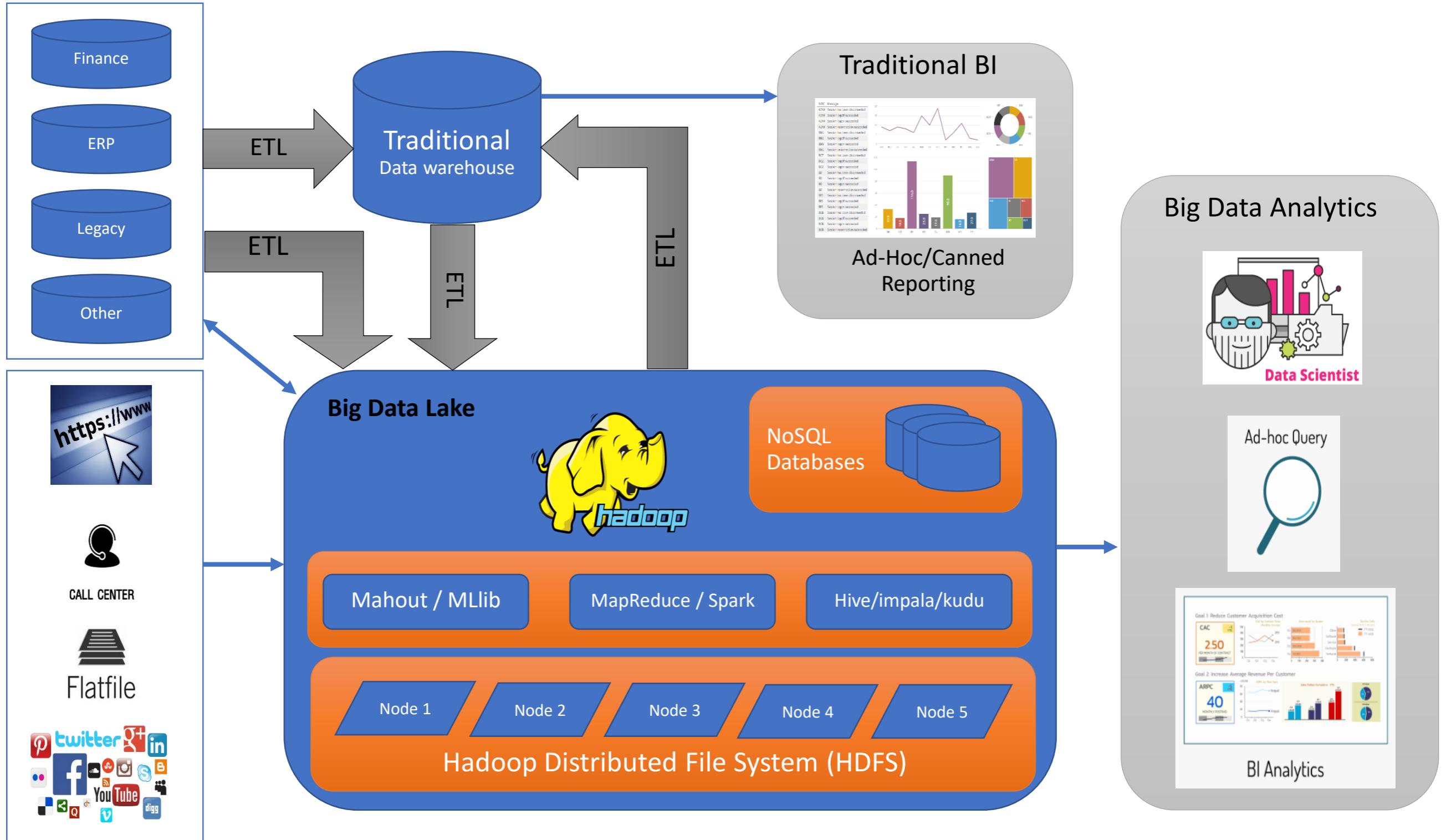
Business Intelligence

- เน้น Descriptive Analytics และ Diagnostic Analytics
- เน้นความรู้สถิติ
- ถ้ากษณะข้อมูล
 - มีจำกัด
 - Cleaned data
 - Structure
 - Simple model
- เทคโนโลยี
 - Data warehouse

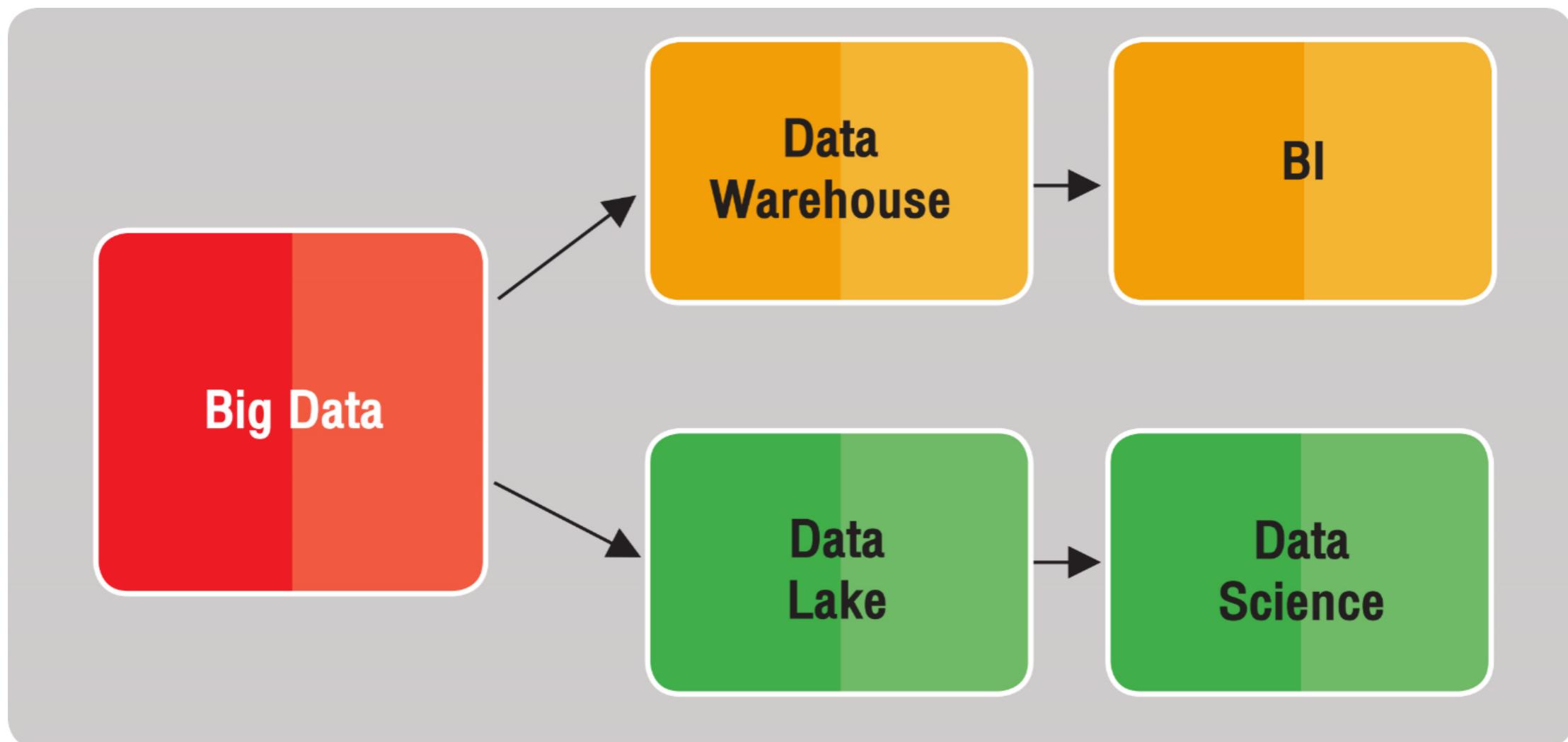
Big Data Analytics

- เน้น Predictive Analytics และ Prescriptive Analytics
- เน้น Data Science
- ถ้ากษณะข้อมูล
 - ขนาดใหญ่
 - ข้อมูลดิบ
 - หลากหลายชนิด
 - Complex model
- เทคโนโลยี
 - Data Lake

Today's business environment requires Big Data



Analytics



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

3 jobs in data

Data Engineer	Data Scientist	Business and Data Analyst
		
หน้าที่ (Responsibilities)		
ออกแบบช่องทางของข้อมูล และวิธีการจัดเก็บ และใช้งาน	ออกแบบโมเดลจากข้อมูล เพื่อหา ช่องทางใหม่ๆ ให้องค์กร	วิเคราะห์ และออกแบบการ นำเสนอข้อมูล เพื่อแก้ปัญหา ในส่วนต่างๆ ของธุรกิจ
การศึกษา (Education)		
- วิศวกรรมศาสตร์คอมพิวเตอร์ - วิทยาศาสตร์คอมพิวเตอร์ - เทคโนโลยีระบบข้อมูล	- วิศวกรรมอุตสาหการ - การวิจัยเชิงปฏิบัติการ - วิทยาศาสตร์คอมพิวเตอร์ - สด็ต	- การตลาด - เศรษฐศาสตร์ - การเงิน - สด็ต - โลจิสติกส์ และอื่นๆ
คุณสมบัติ (Skills)		
- Database design - Production coding - Data warehousing - Data transformation - Cloud computing - Big data design	- Mathematics - Propability - Statistics - Machine Learning - Programming - Data Visualization - Business	- Analysis - Statistics - Business - Communication - Project management
เครื่องมือ (Program)		
- Oracle - Hadoop - SQL - NoSQL - Pig - Hive - Spark - MongoDB - SAP - Pentaho - Cloudera - JavaScript	- Python - R - SAS - MATLAB - Julia - Spark - Tableau	- R - SPSS - SAS - Tableau - Microsoft Excel