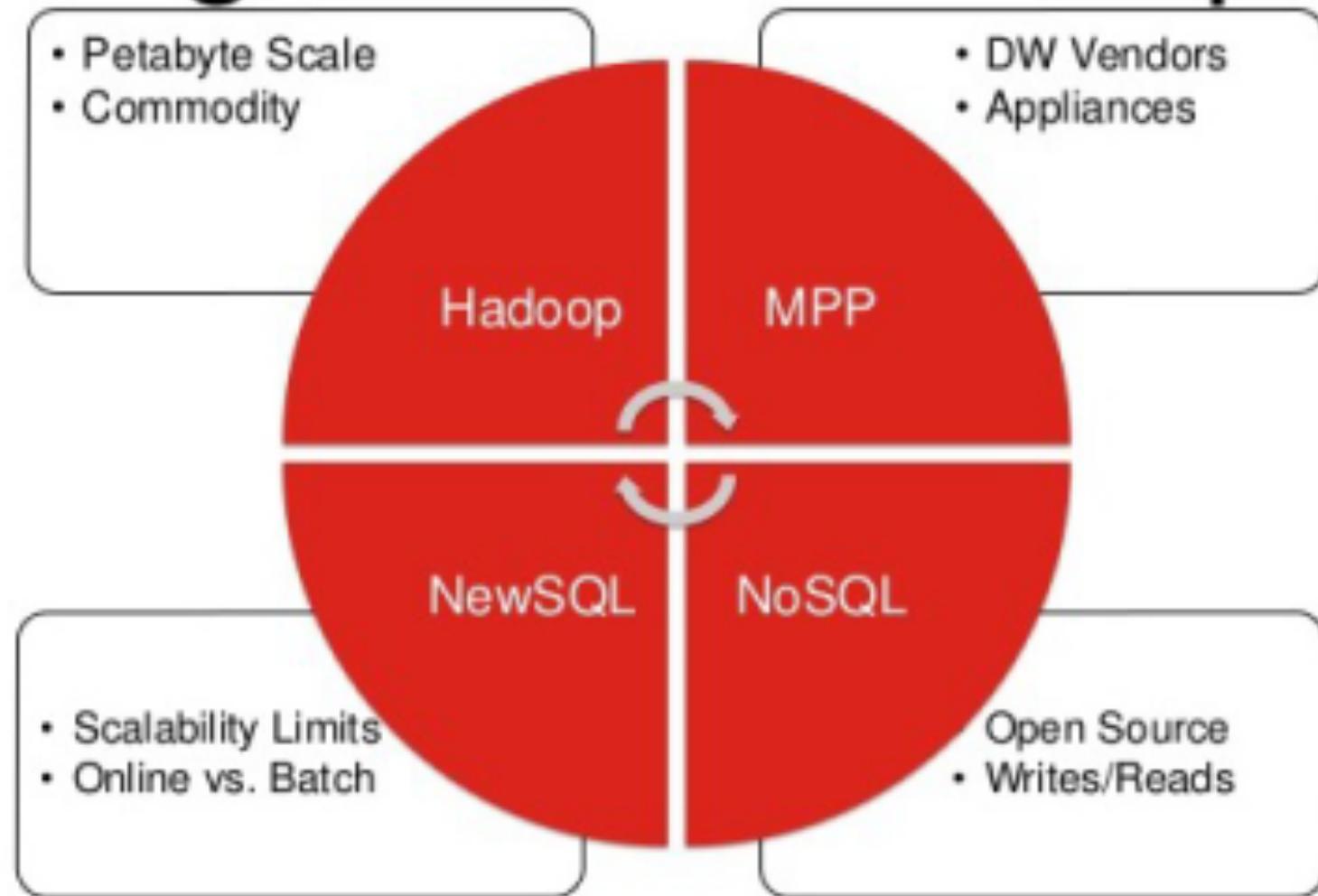


Big Data Technology

Big Data Landscape



Big Data Landscape 2017

BIG DATA LANDSCAPE 2017



CHULALONGKORN
BUSINESS SCHOOL



Last updated 4/5/2017

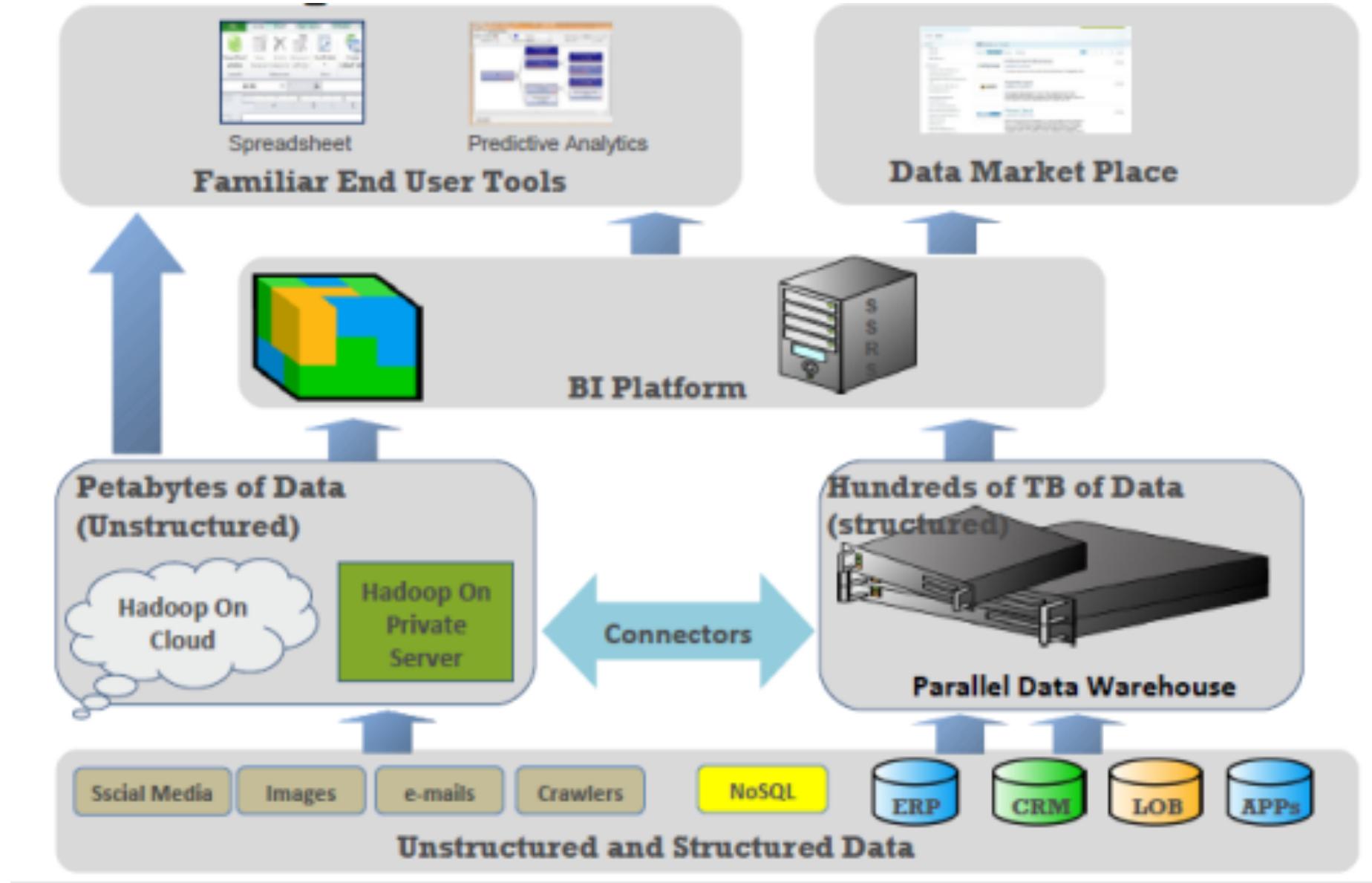
© Matt Turck (@mattturck), Jim Hao (@jimrhao), & FirstMark (@firstmarkcap)

mattturck.com/bigdata2017

FIRSTMARK
EARLY STAGE VENTURE CAPITAL

Thaveewat Khanan

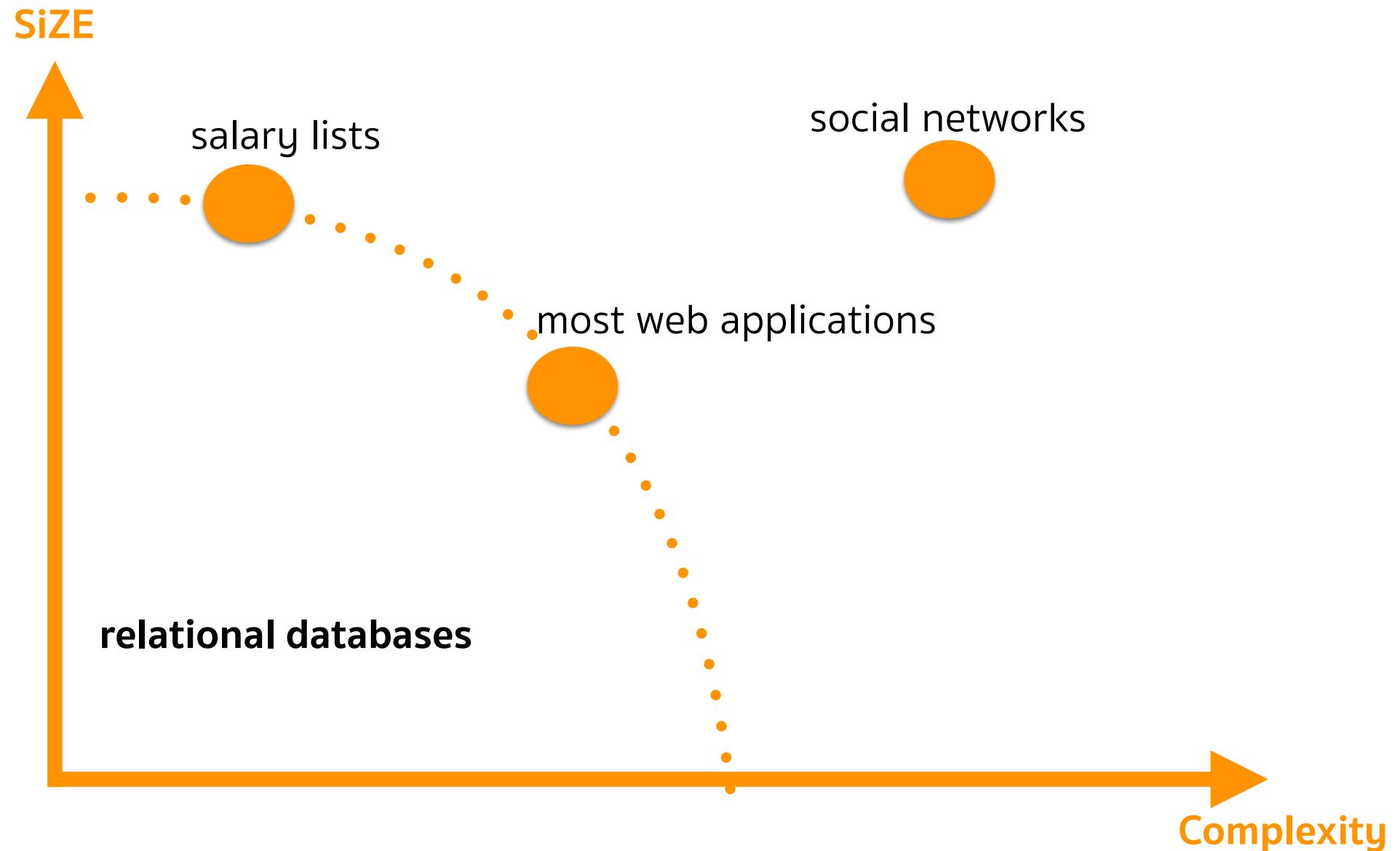
Big Data Future Architecture



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)

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
1	Ricky
2	Ankur
3	Sam

id	Age
2	20
3	25

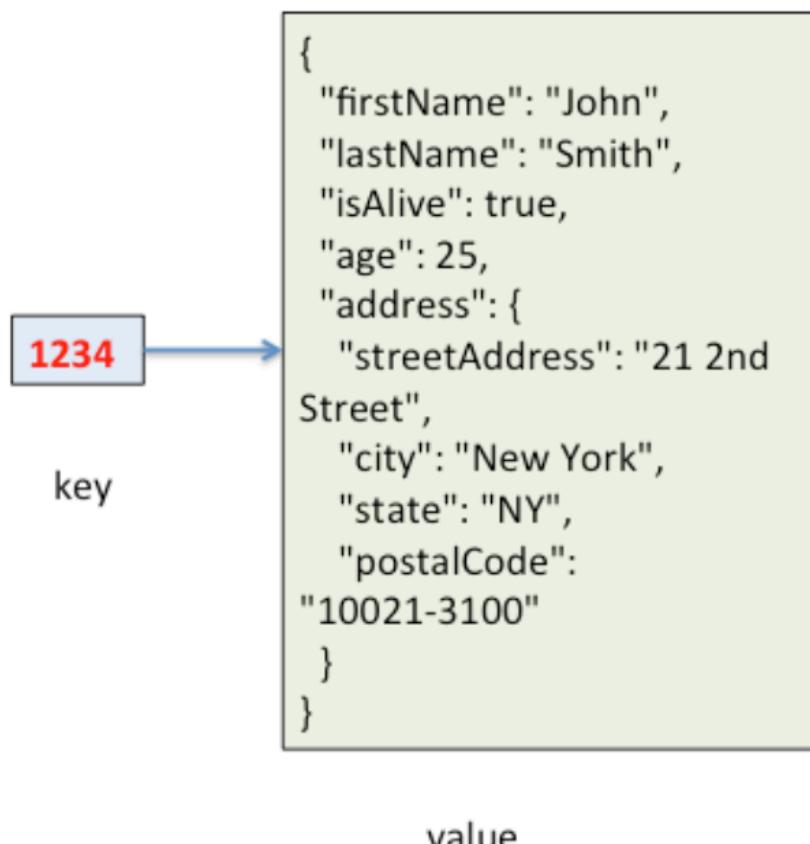
id	Interests
1	Soccer
1	Movies
1	Baseball
3	Music



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": "SML",  
    "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



Author name

Comment Table
Reader Table

Article Table
Author Table

Relational Database approach
Document store approach

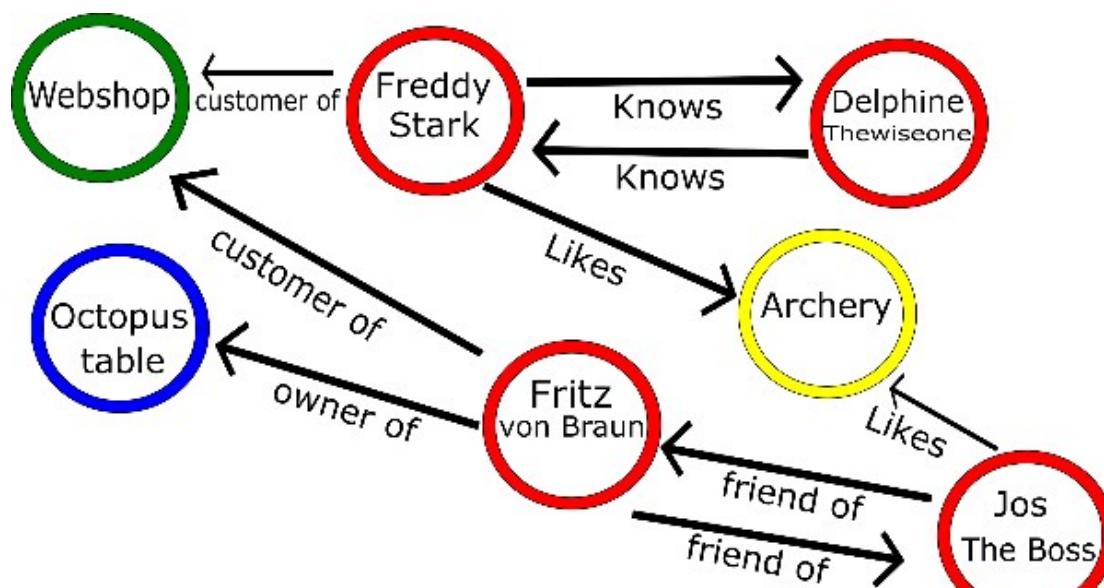
Whereas relational databases chop up data, Document stores save documents as a single entity

```
{  
  "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"  
        }  
      ]  
    }  
  ]  
}
```



Graph database

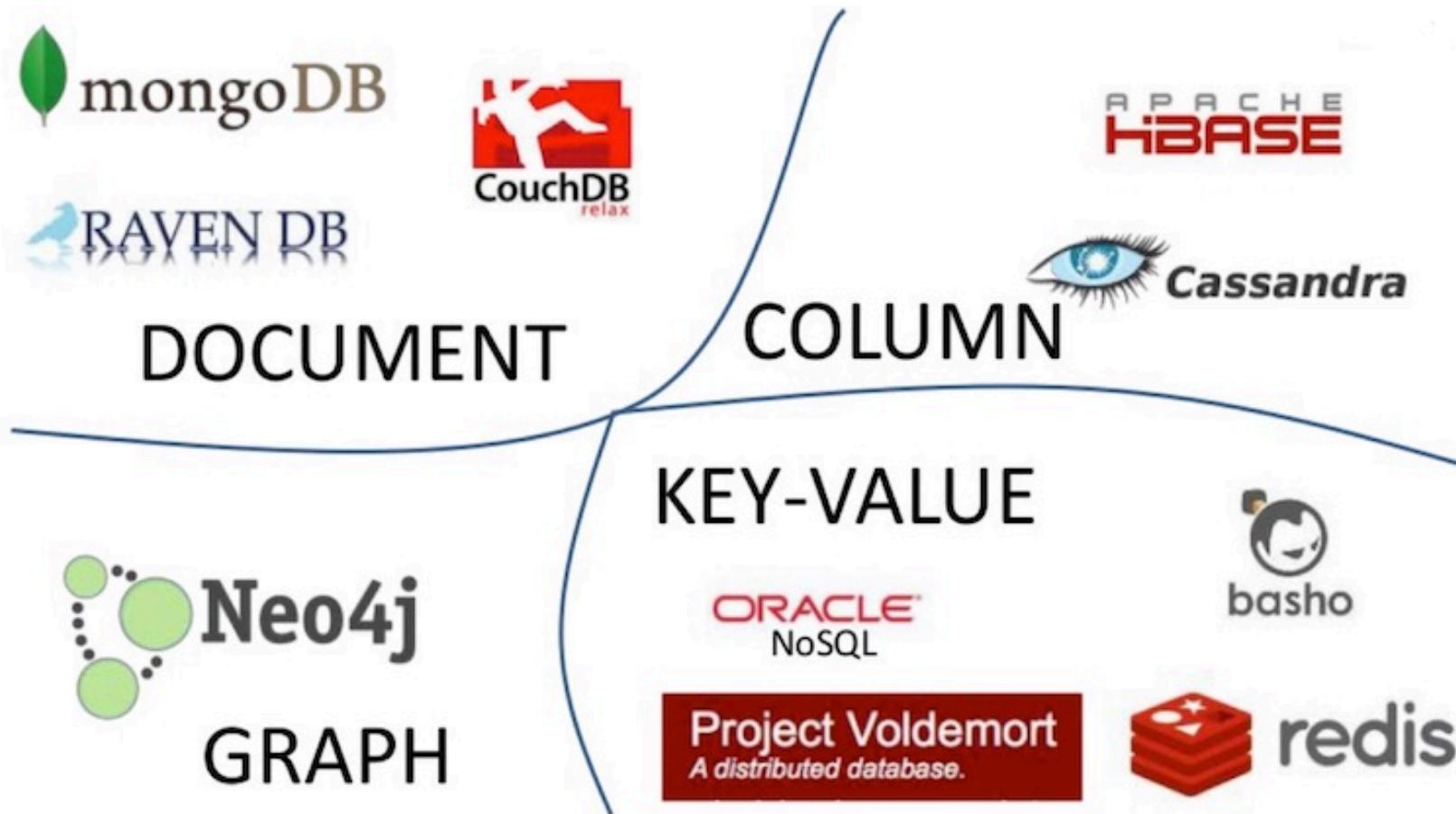
A database that uses graph structures for semantic queries with nodes, edges, and properties to represent and store data.

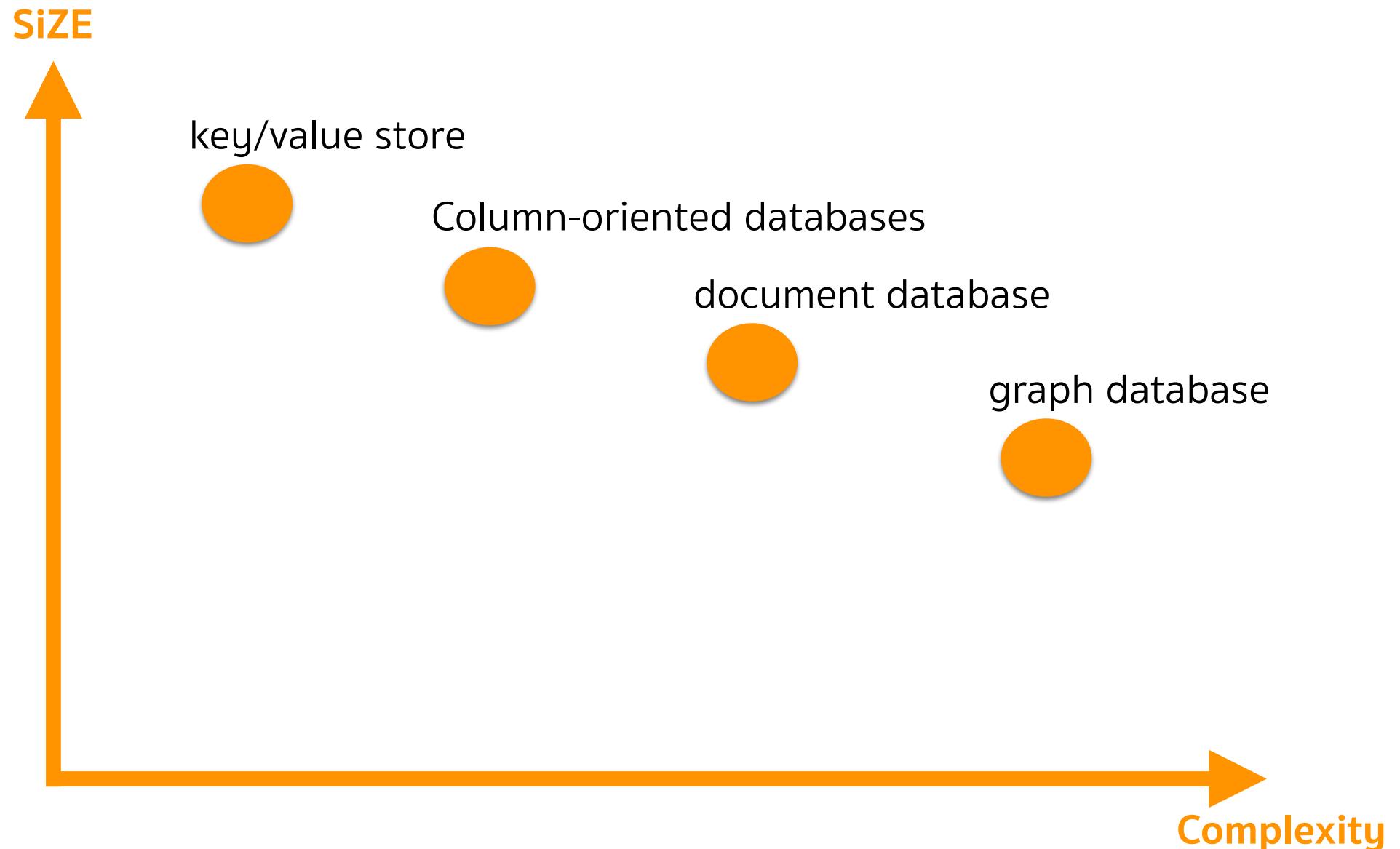


 Neo4j
the graph database

 AllegroGraph
Franz Inc.

 InfiniteGraph





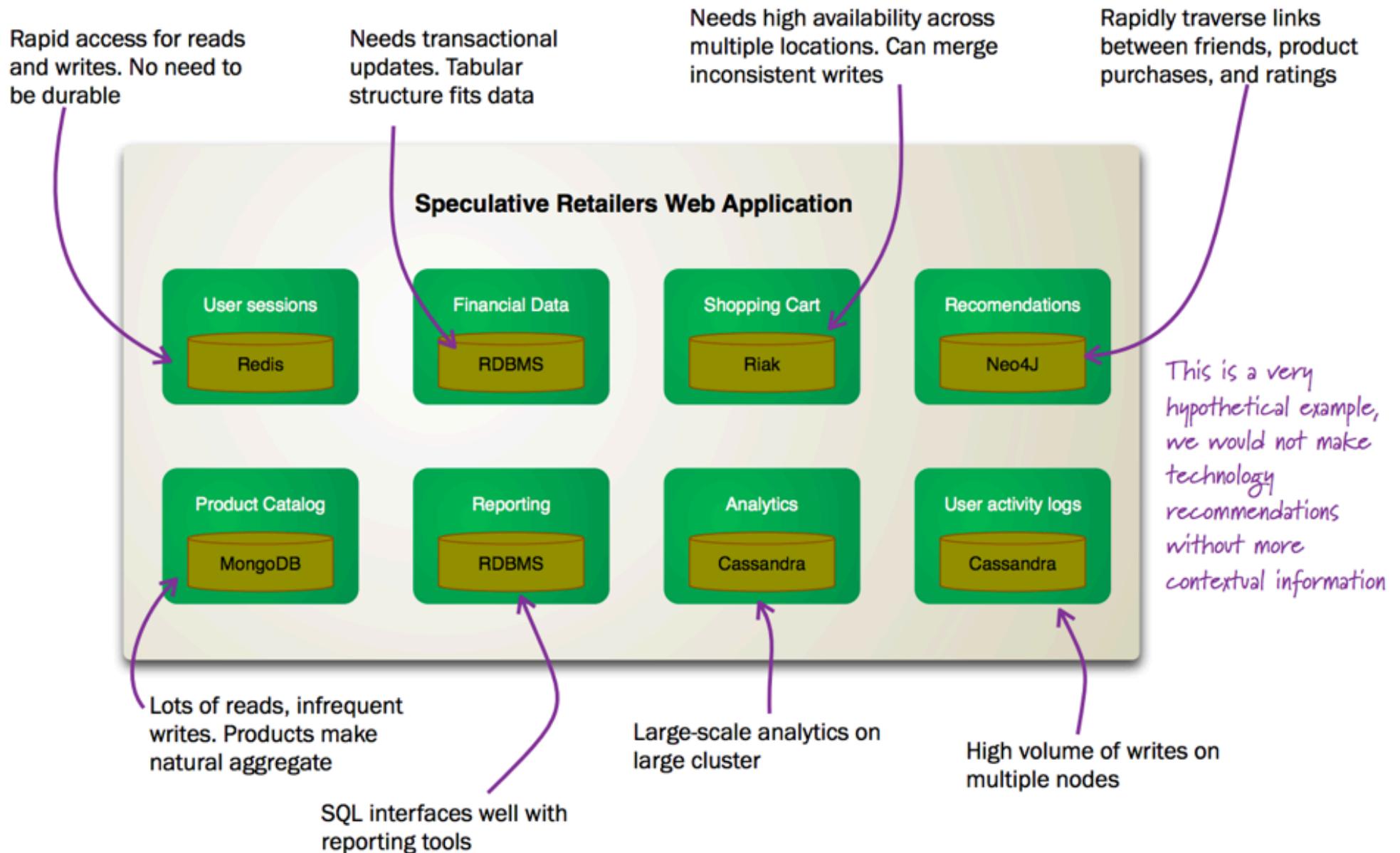
SQL

works great, can't scale for large data

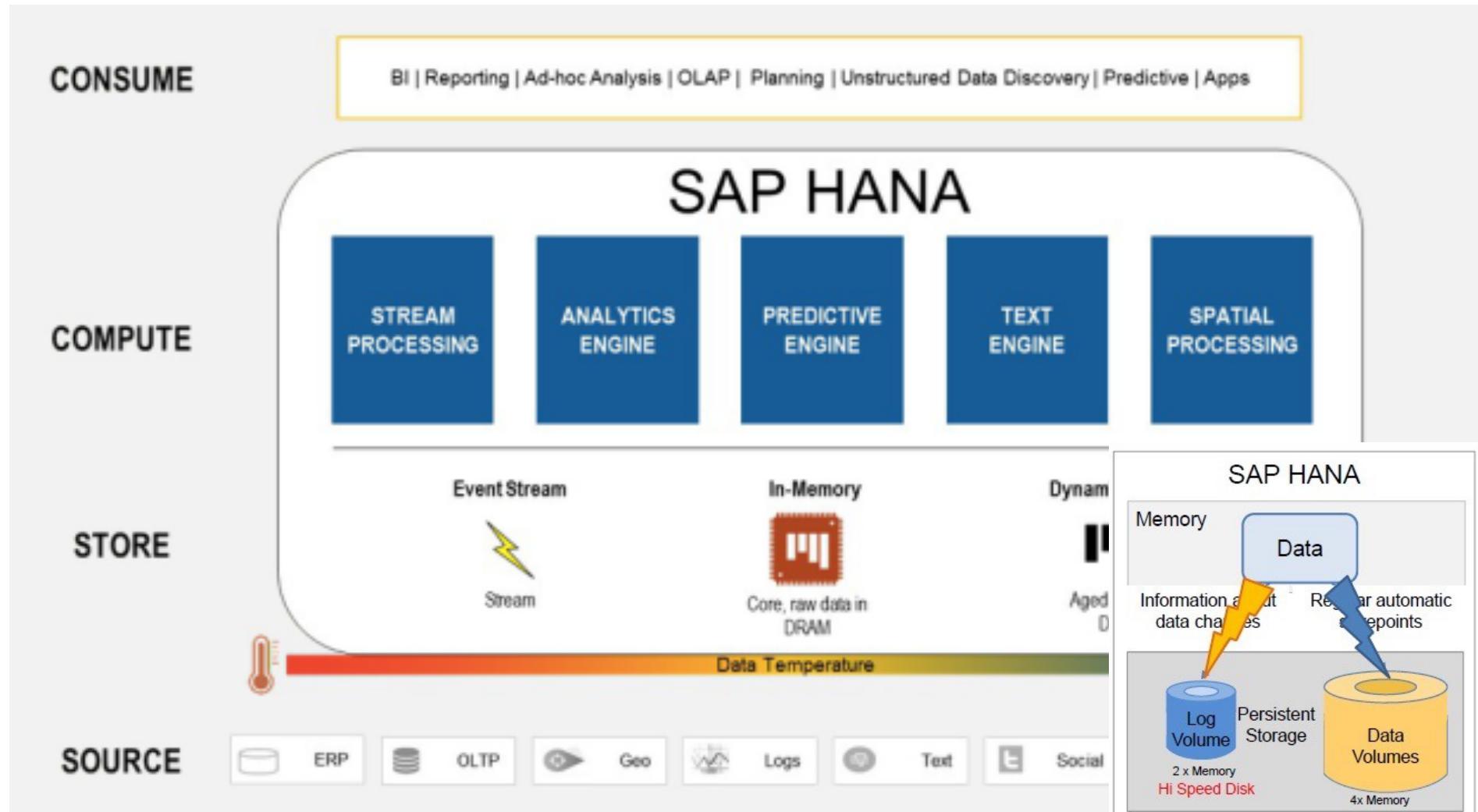
NoSQL

works great, doesn't fit all situations

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



SAP Hana



MPP:Oracle Exadata

Oracle Exadata Database Machine

Extreme Performance for the Cloud

Ellison announces next-generation systems



Ease compliance: OFSAA and Oracle Exadata (PDF)



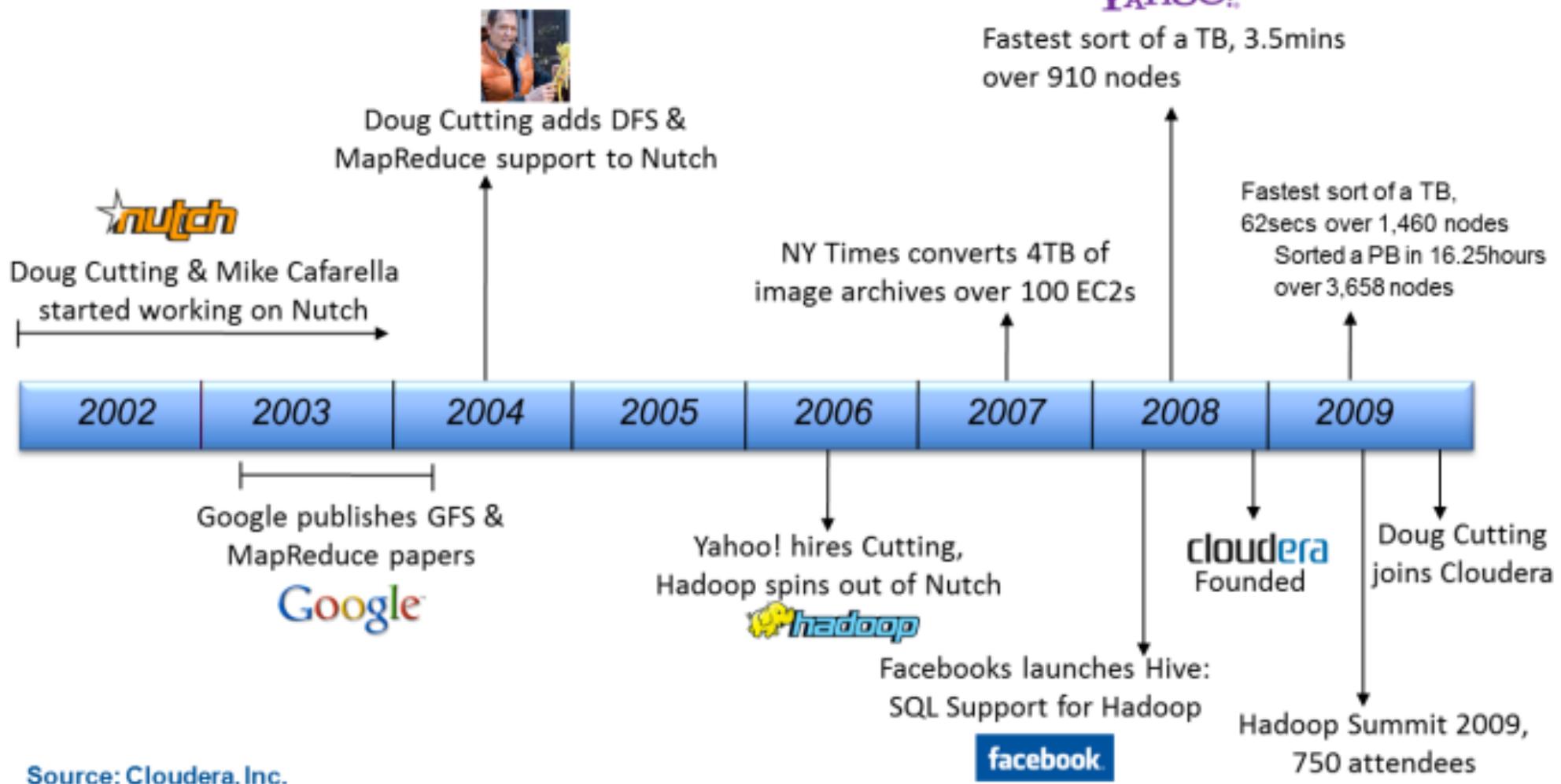
What is Hadoop ?

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

Completely written in java
Open source & distributed under Apache license

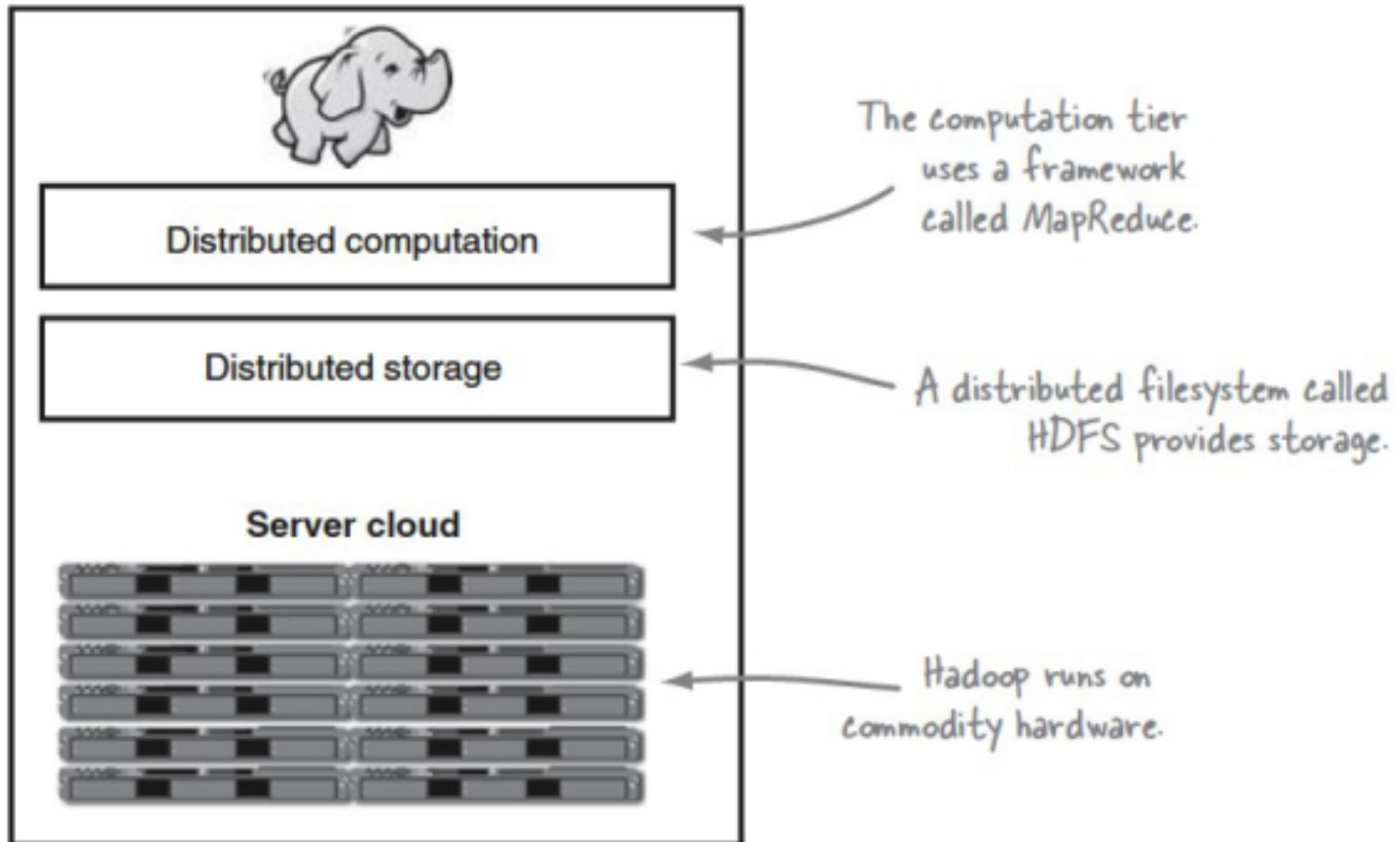


Hadoop Creation History

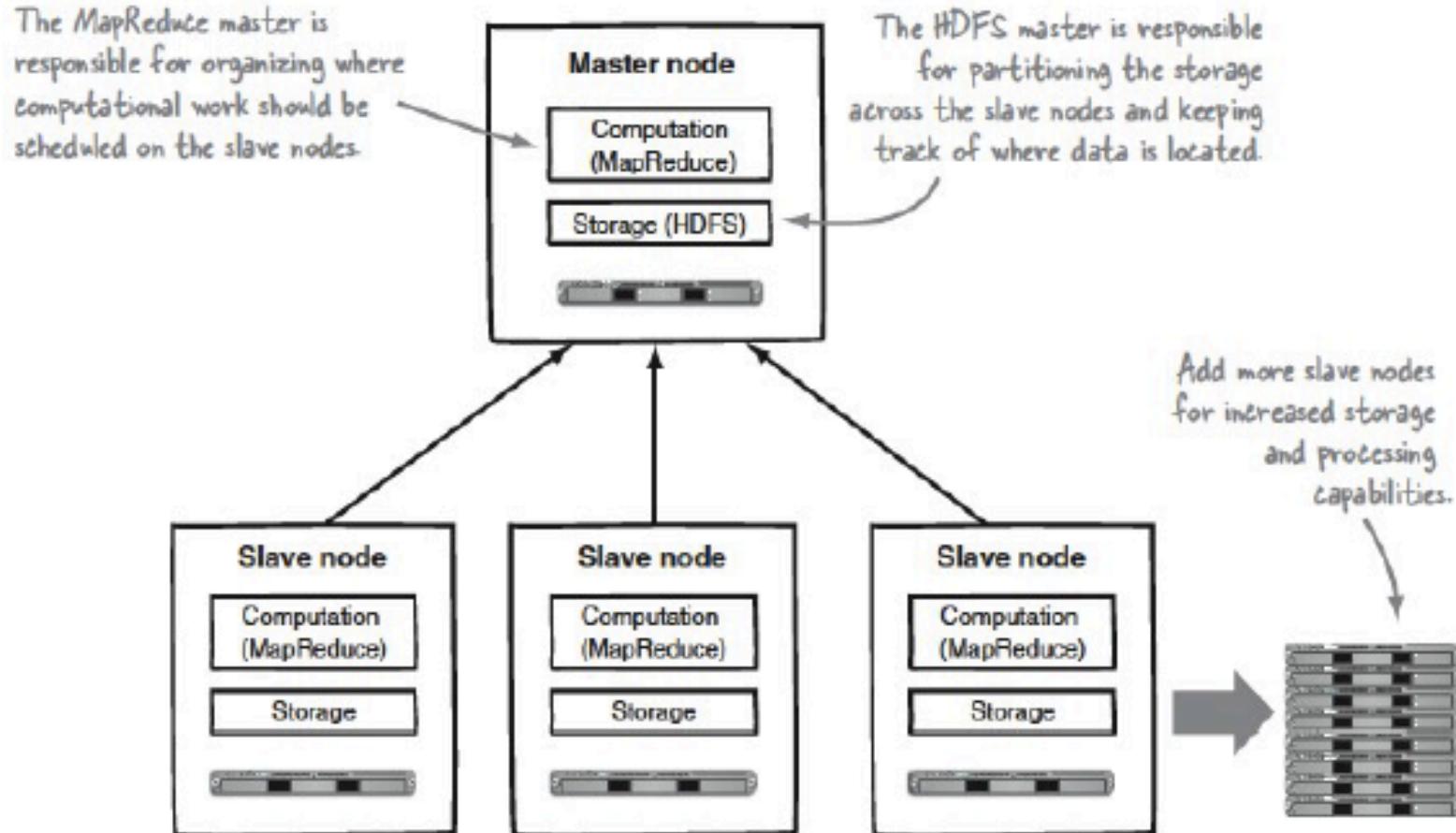


Source: Cloudera, Inc.

Hadoop Environment



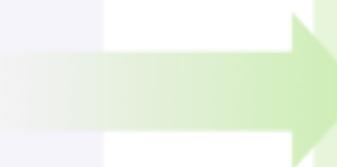
Hadoop Architecture



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

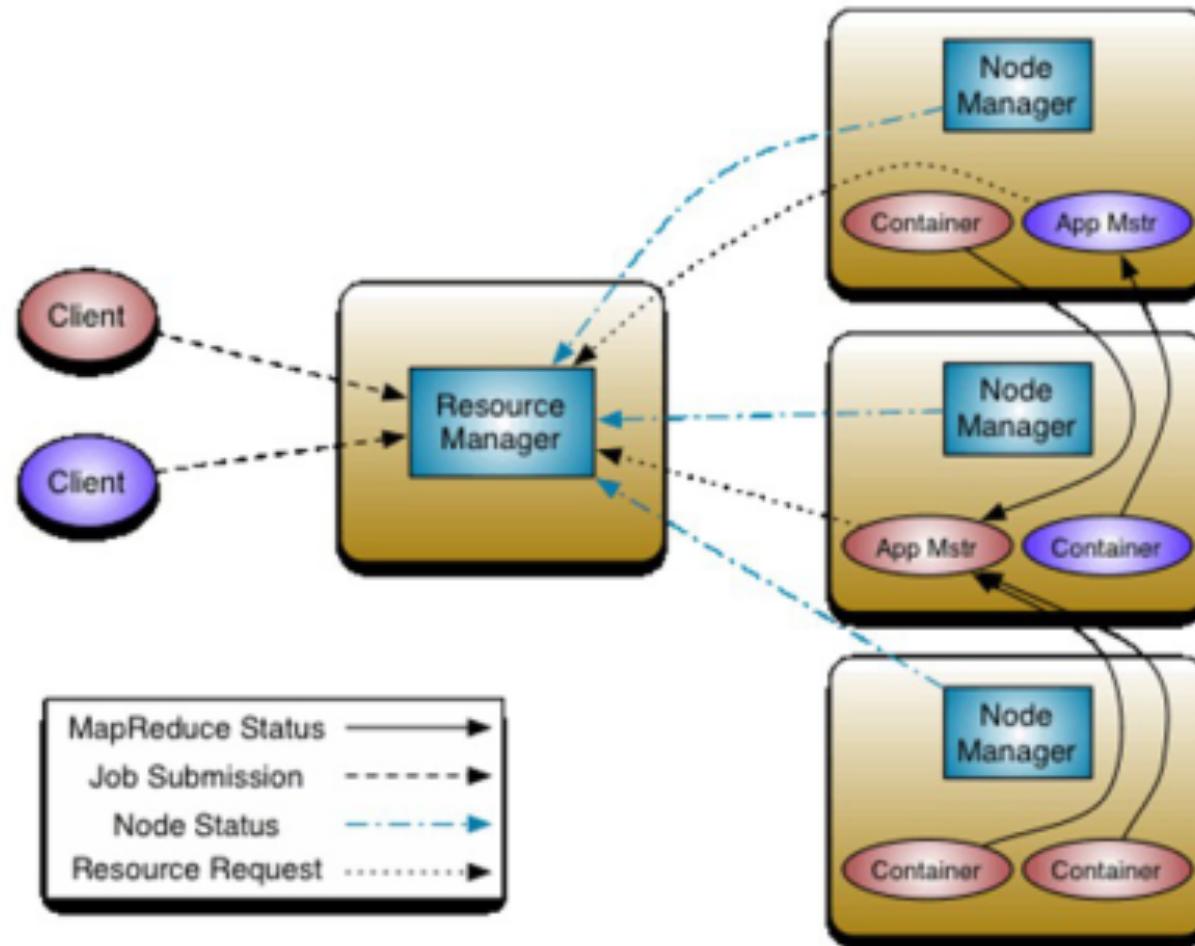
YARN: Data Operating System
(Cluster Resource Management)

MapReduce
(Cluster Resource Management
& Batch Data Processing)

1
HDFS
(Hadoop Distributed File System)

1
HDFS
(Hadoop Distributed File System)
N

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

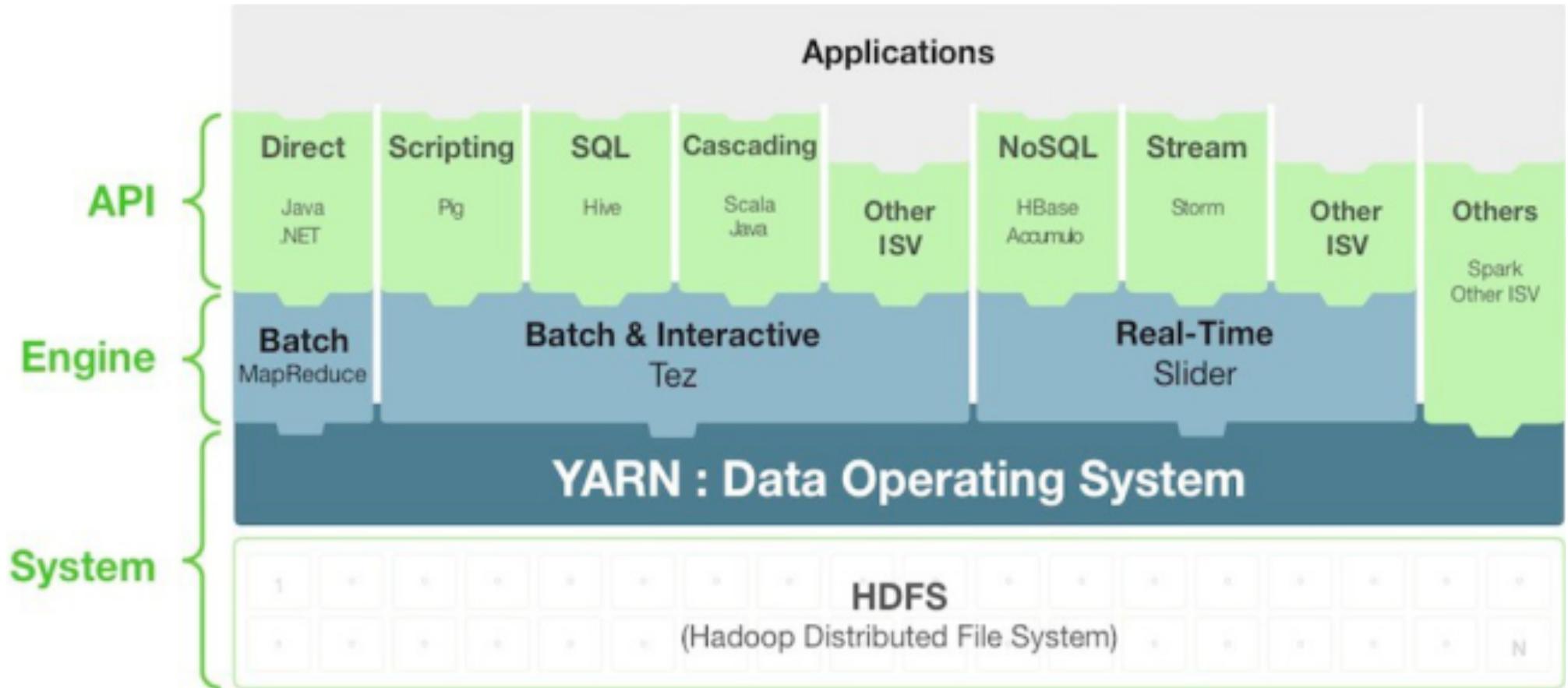
The stack is continually evolving and growing!

2006	2007	2008	2009	2010	2011	2012	2013	2014-15
Core Hadoop (HDFS, MapReduce)	Solr Pig	Core Hadoop	HBase ZooKeeper	Hive Mahout	Hive Mahout	Mahout	Mahout	Mahout
			Solr Pig	Solr Pig	Solr Pig	HBase	HBase	HBase
			ZooKeeper	ZooKeeper	ZooKeeper	ZooKeeper	ZooKeeper	ZooKeeper
					Solr	Solr	Solr	Solr
				Pig	Pig	Pig	Pig	Pig
				Pig	YARN	YARN	YARN	YARN
				Core Hadoop	Core Hadoop	Core Hadoop	Core Hadoop	Core Hadoop
					Flume	Flume	Flume	Flume
					Bigtop	Bigtop	Bigtop	Bigtop
					Oozie	Oozie	Oozie	Oozie
					MRUnit	MRUnit	MRUnit	MRUnit
					HCatalog	HCatalog	HCatalog	HCatalog
					Hue	Hue	Hue	Hue
					Sqoop	Sqoop	Sqoop	Sqoop
					Whirr	Whirr	Whirr	Whirr
					Avro	Avro	Avro	Avro
					Hive	Hive	Hive	Hive
					Mahout	Mahout	Mahout	Mahout
					HBase	HBase	HBase	HBase
					ZooKeeper	ZooKeeper	ZooKeeper	ZooKeeper
					Solr	Solr	Solr	Solr
					Pig	Pig	Pig	Pig
					YARN	YARN	YARN	YARN
					Core Hadoop	Core Hadoop	Core Hadoop	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	

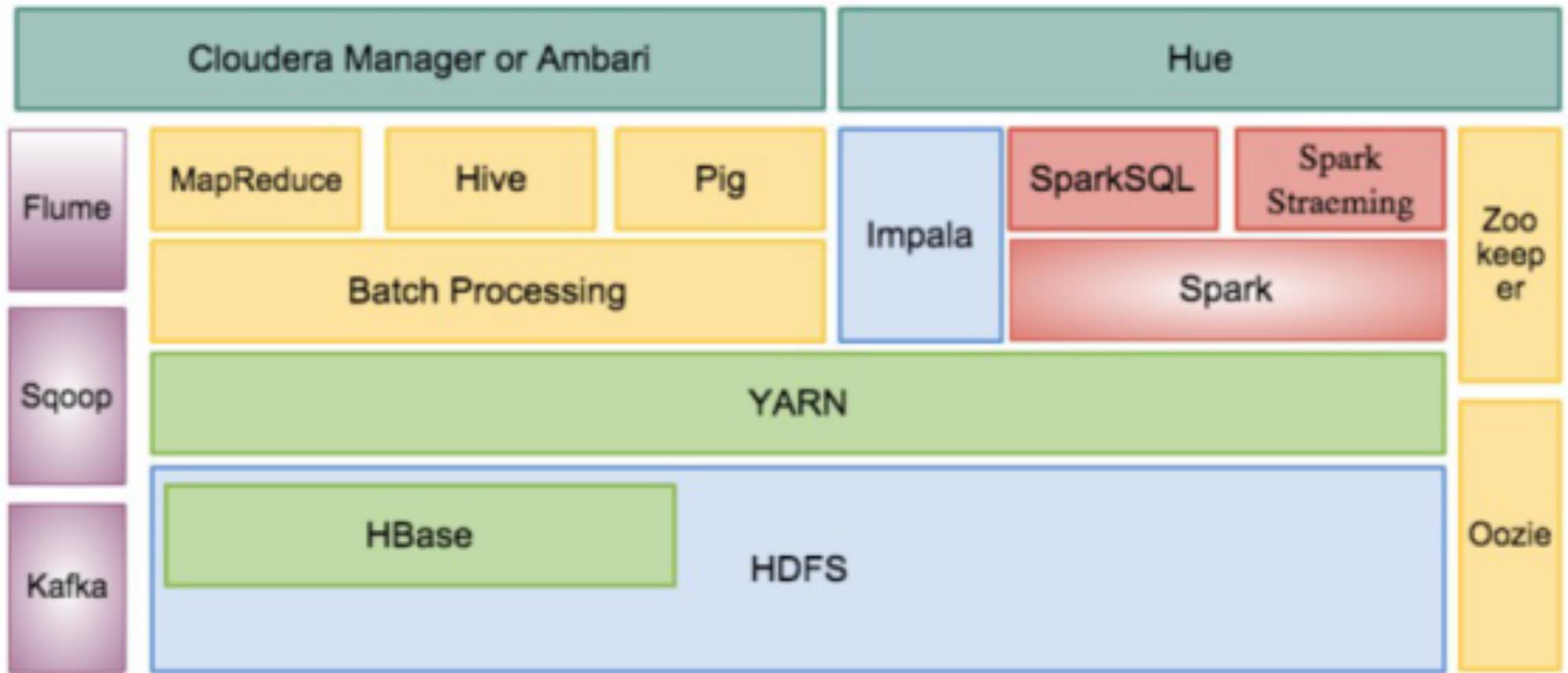
cloudera

© Cloudera, Inc. All rights reserved. 9

Hadoop 2.x Ecosystems



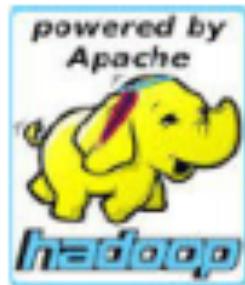
Hadoop Ecosystems



Hadoop Distribution



CHULALONGKORN
BUSINESS SCHOOL
FLAGSHIP FOR LIFE



MAPR

cloudera



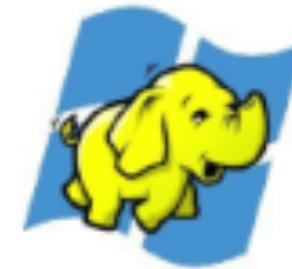
Pivotal™



TERADATA

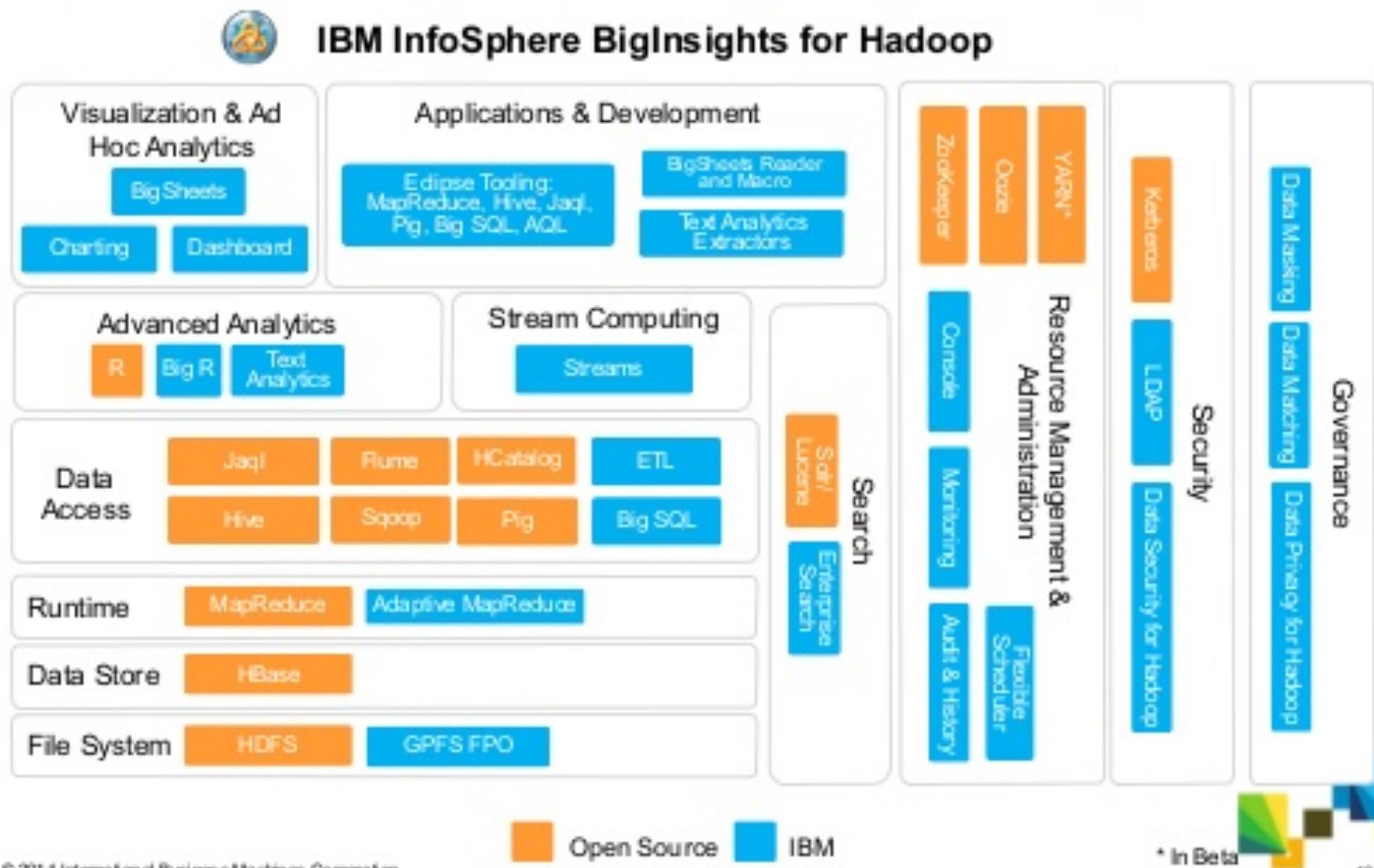


amazon
web services™

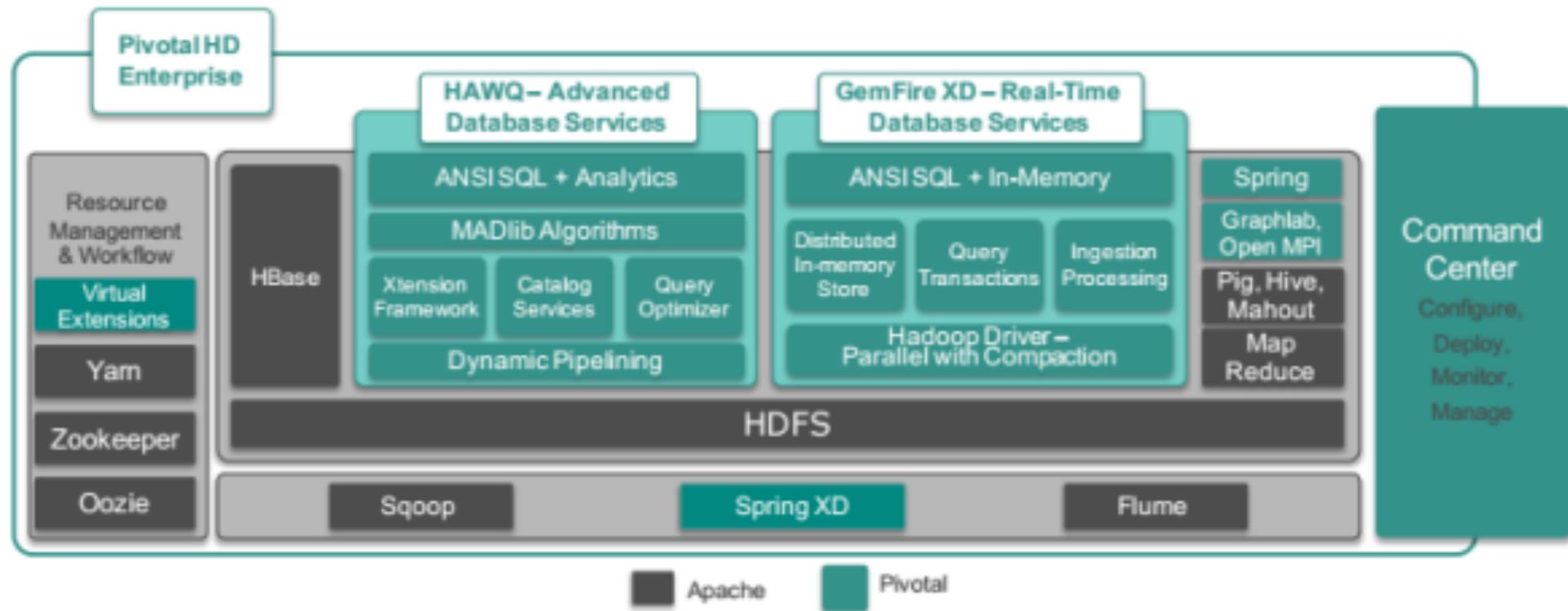


Microsoft Azure

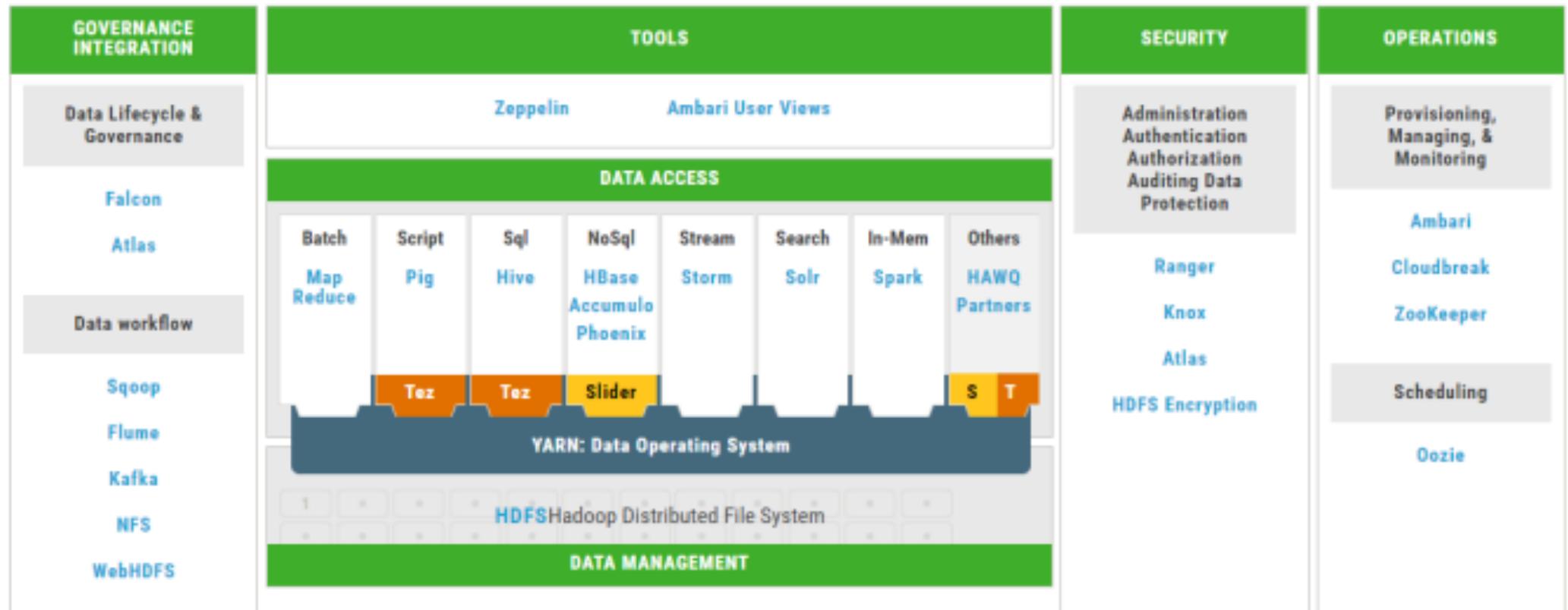
IBM InfoSphere BigInsights



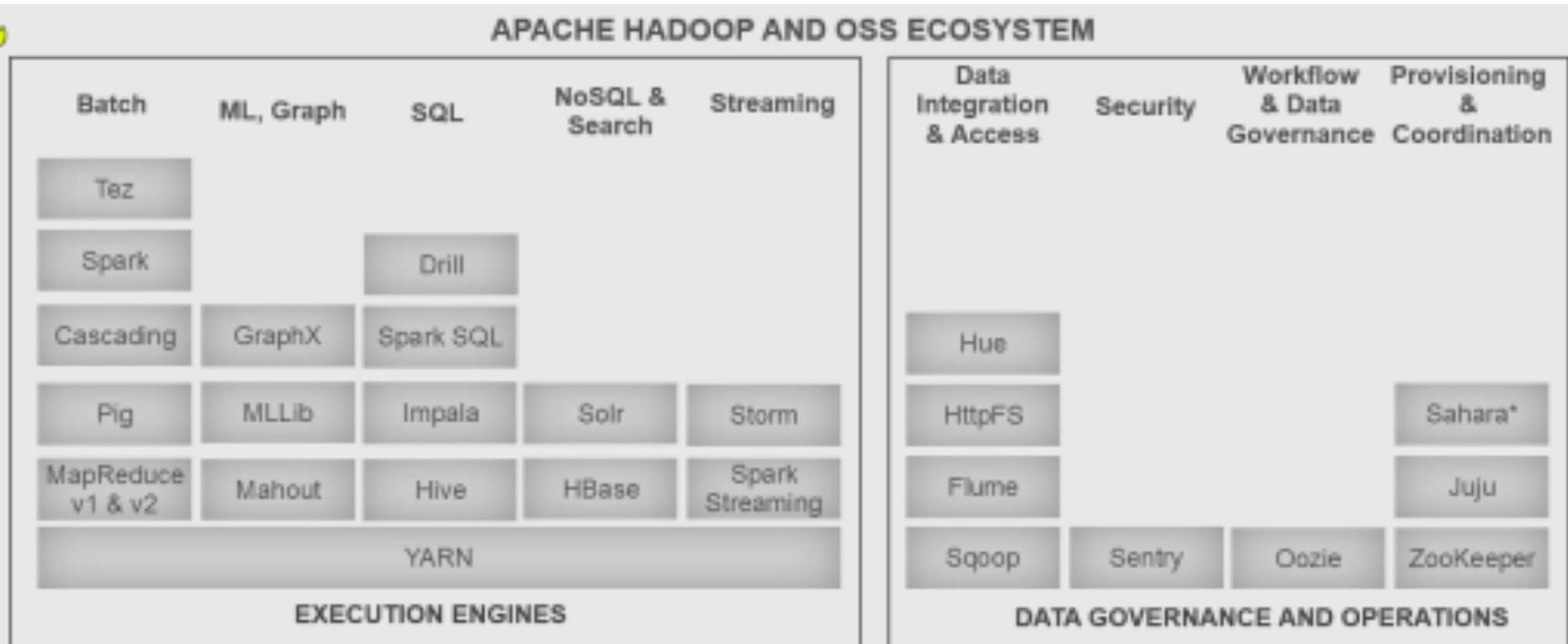
Pivotal HD Architecture



Hortonworks



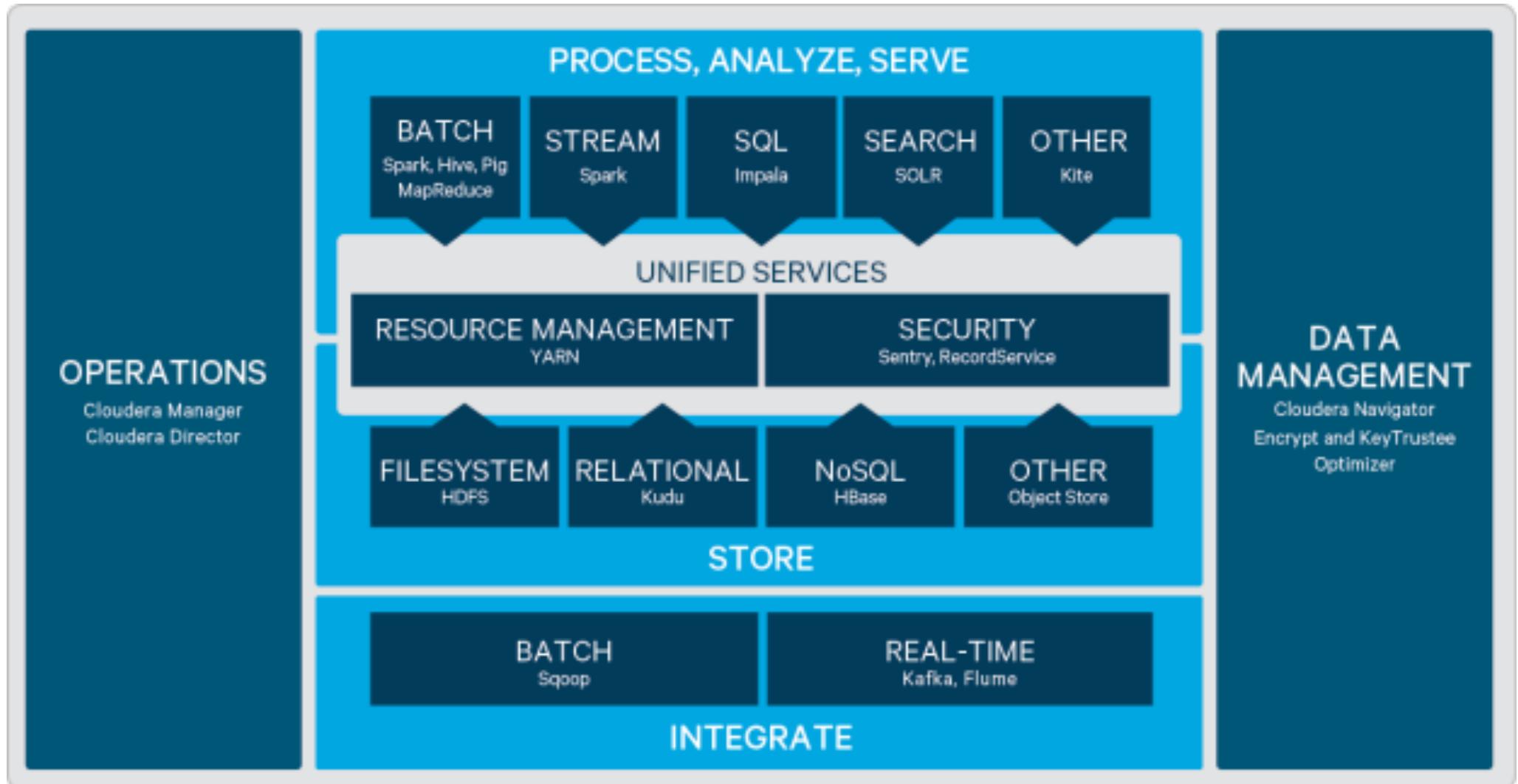
Management



MapR-FS

Data Platform

MapR-DB



Default Cloudera Services

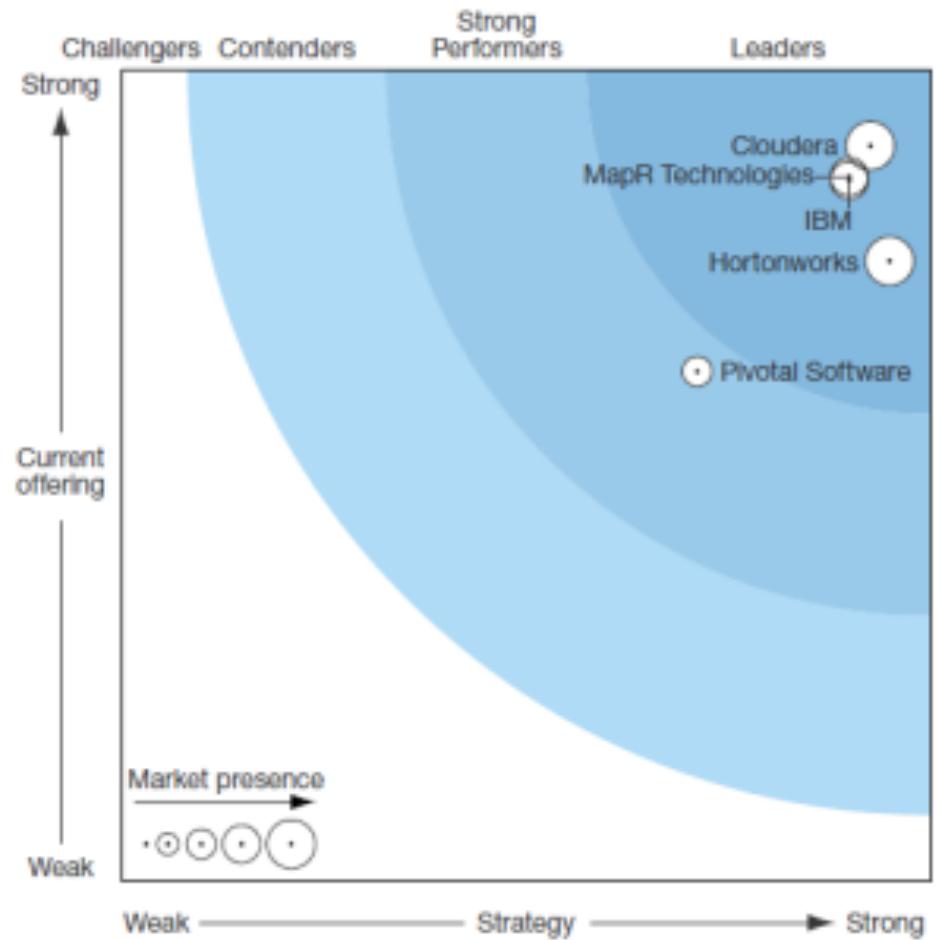
- Cloudera Manager
- HDFS
- YARN
- Apache Hive and Pig
- Apache Flume and Sqoop
- Apache Oozie
- Cloudera Hue
- ZooKeeper



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

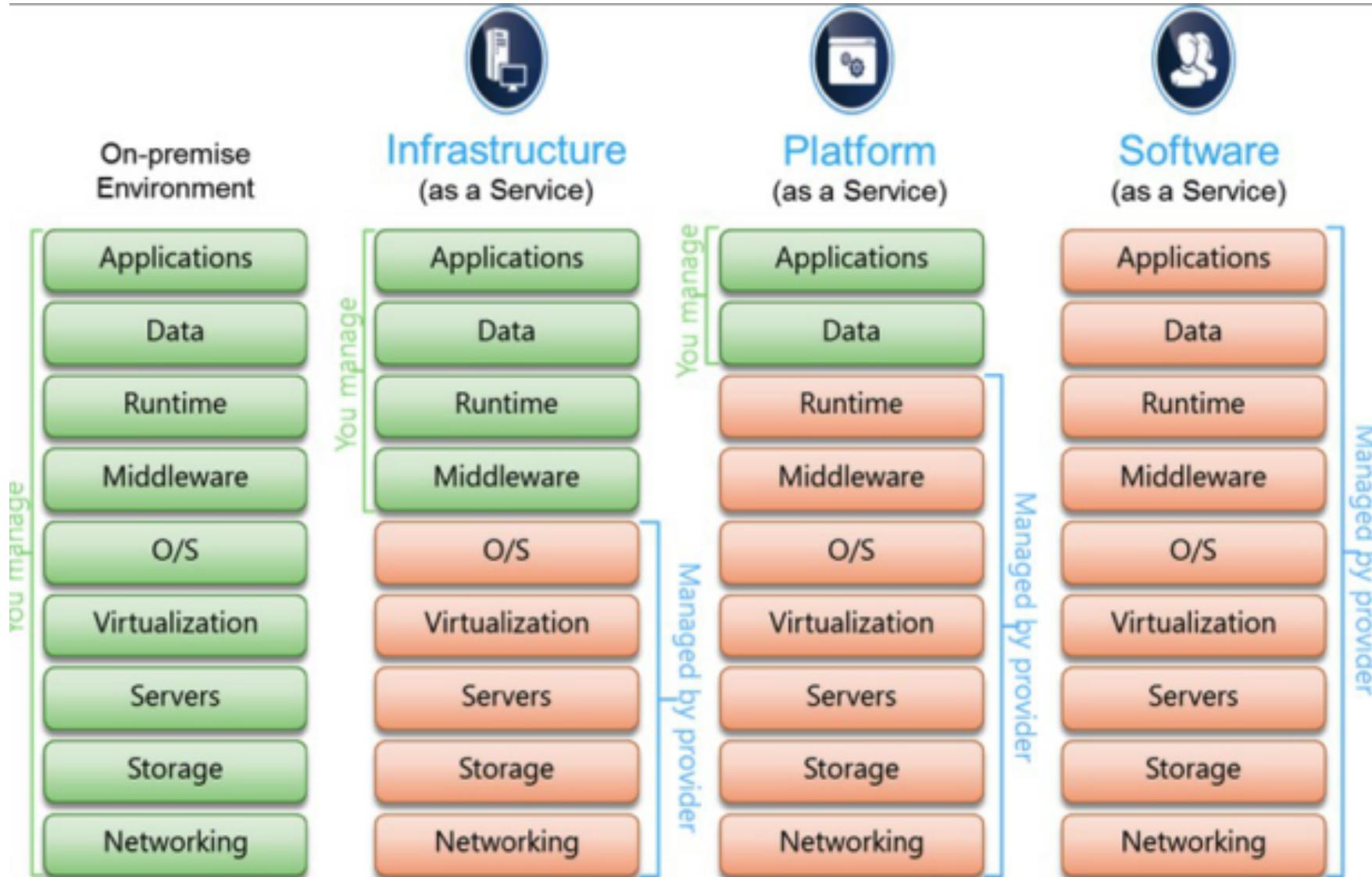


Big Data Using Public Cloud

Issue with Big Data Infrastructure

- Large investment
- Scalability
- ROI
- Business Cases

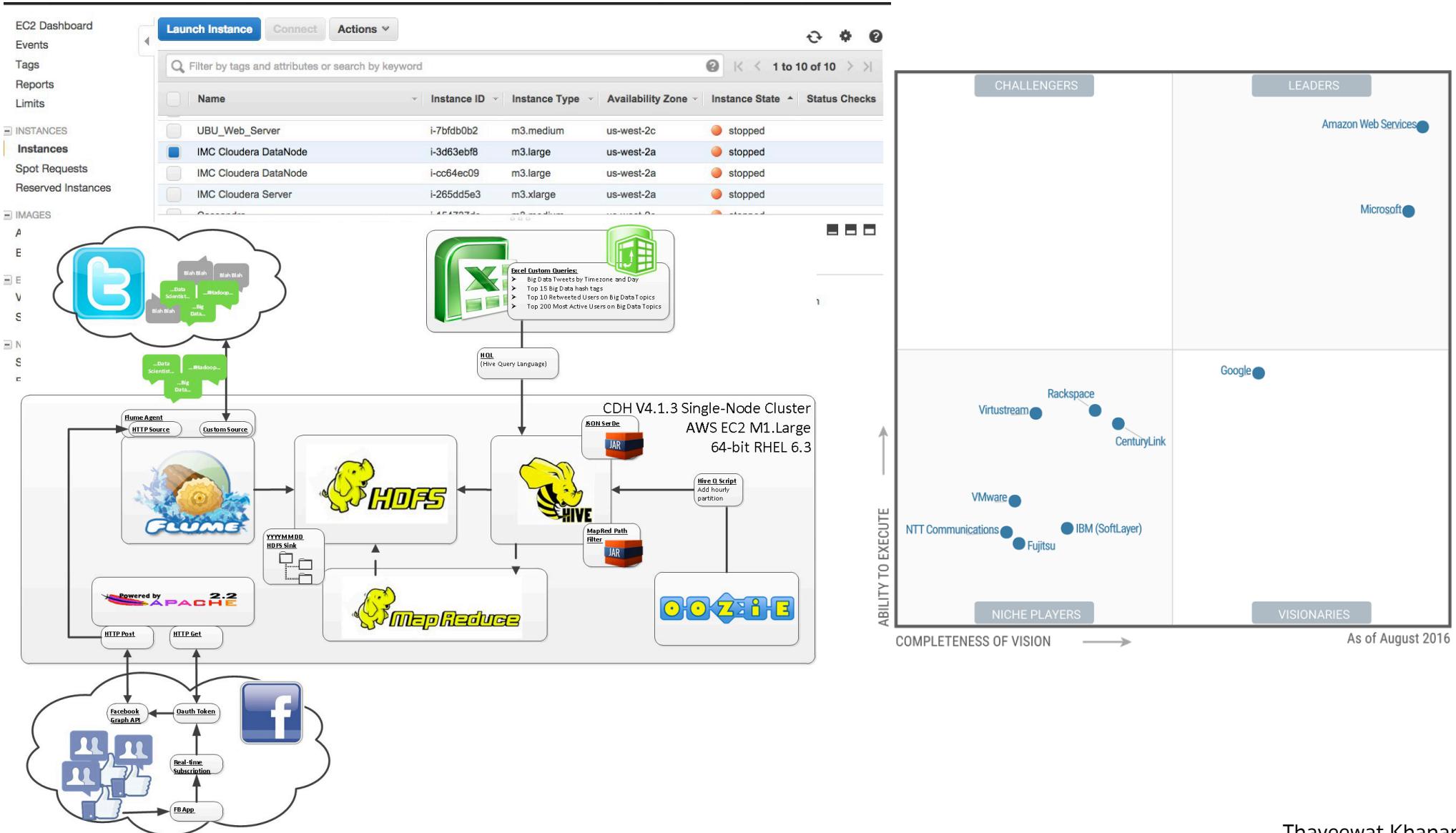
Cloud Technology





Big Data on Cloud

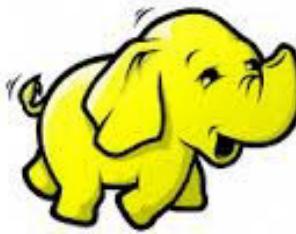
Big Data using IaaS



Big Data on Cloud

Using Big Data as a Services

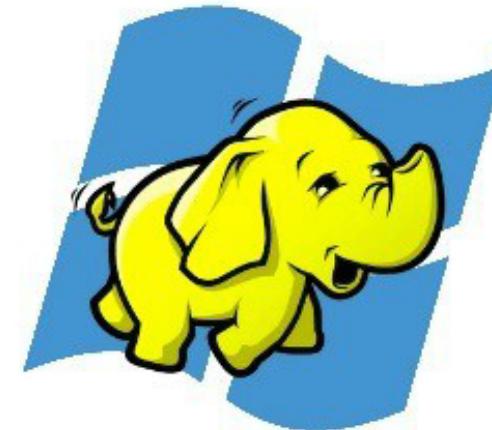
Amazon Web Services



Amazon

Elastic Mapreduce

Google bigquery



Microsoft Azure Hadoop