

# **Chapter3**

# **Understanding Big Data Technology Foundations**

## **Big Data Stack**

**Basanta Joshi, PhD**

Asst. Prof., Depart of Electronics and Computer Engineering  
Program Coordinator, MSc in Information and Communication Engineering  
Member, Laboratory for ICT Research and Development (LICT)

Member, Research Management Cell (RMC)

Institute of Engineering

[basanta@ioe.edu.np](mailto:basanta@ioe.edu.np)

<http://www.basantajoshi.com.np>

<https://scholar.google.com/citations?user=iocLiGcAAAAJ>

[https://www.researchgate.net/profile/Basanta\\_Joshi2](https://www.researchgate.net/profile/Basanta_Joshi2)

# Big data Stack

VISUALIZATION  
LAYER AND API'S



RESTful API

ANALYTICAL TOOLS



ANALYSIS LAYER



PROCESSING  
LAYER



STORAGE LAYER



REAL TIME  
EVENT ENGINE



INTEGRATION LAYER  
(ETL)



SOURCES



Logs XML RSS Twitter Facebook Google+

CAPA DE GESTIÓN



HERRAMIENTAS  
ADMINISTRACIÓN



HERRAMIENTAS  
MONITORIZACIÓN



Nagios

HERRAMIENTAS  
DIAGNÓSTICO



# What We Need?

- Store
- Join
- Index
- Analytics
- Aggregate
- Visualize

# Challenge

The challenge in big data analytics is to

- dig deeply
- quickly (real time?)
- and widely

# "ilities" or NFR?

- Availability
- Scalability
- Security
- Performance
- ...

# Solution?

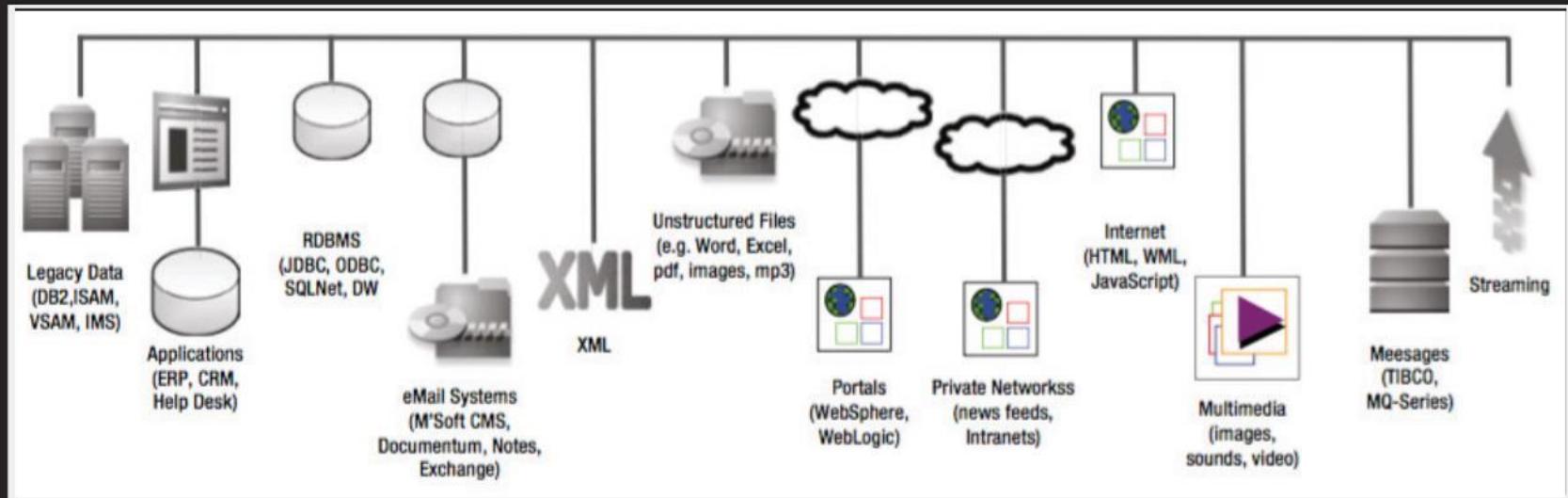
# Big Data Tech Stack

# What're essential components?

# Data Sources

Multiple internal  
& external  
data sources

Creates a  
data lake



## **Legacy Data Sources**

**HTTP/HTTPS web services**

**RDBMS**

**FTP**

**JMS/MQ based services**

**Text/flat file/csv logs**

**XML data sources**

**IM Protocol requests**

## **New Age Data Sources**

### **High Volume Sources**

1. Switching devices data
2. Access point data messages
3. Call data record due to exponential growth in user base
4. Feeds from social networking sites

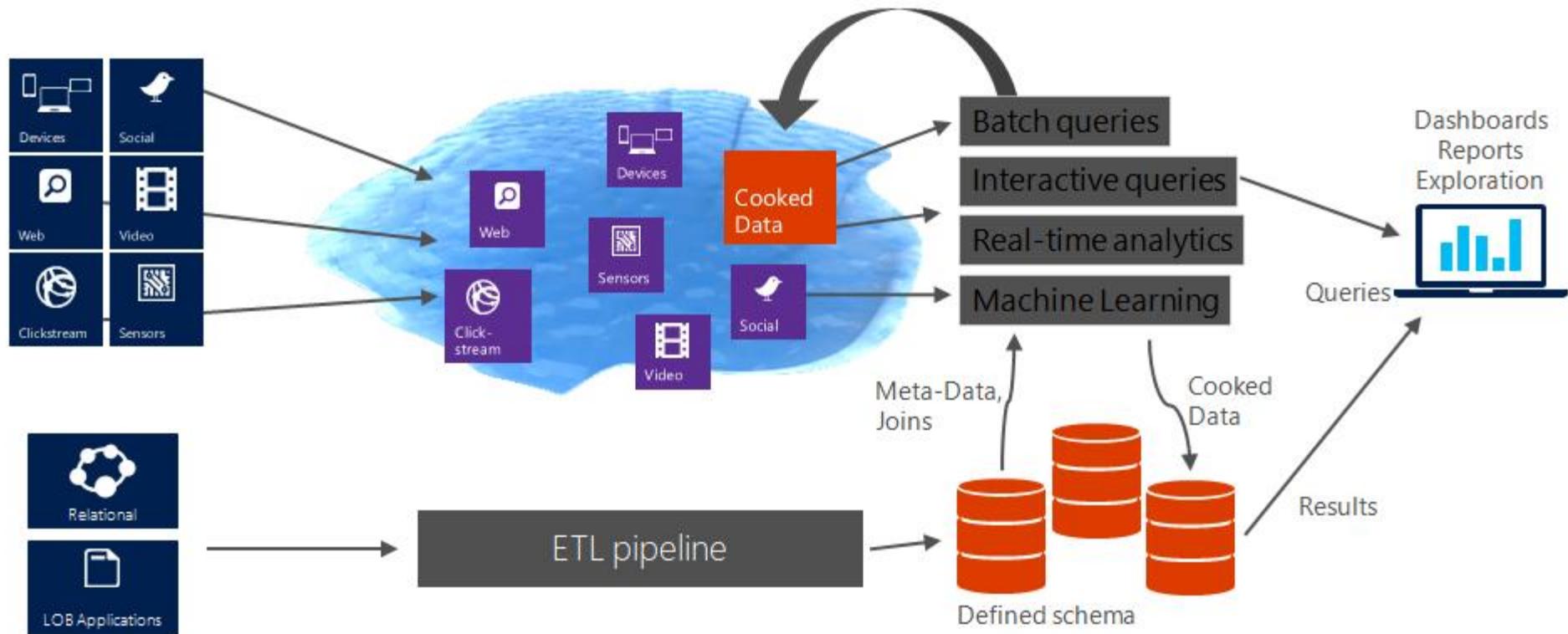
### **Variety of Sources**

1. Image and video feeds from social Networking sites
2. Transaction data
3. GPS data
4. Call center voice feeds
5. E-mail
6. SMS

### **High Velocity Sources**

1. Call data records
2. Social networking site conversations
3. GPS data
4. Call center - voice-to-text feeds

# The data lake and warehouse



Different  
Volume, Variety,  
Velocity

Aim is to create  
a **funnel** after  
proper **validation**  
and **cleaning**

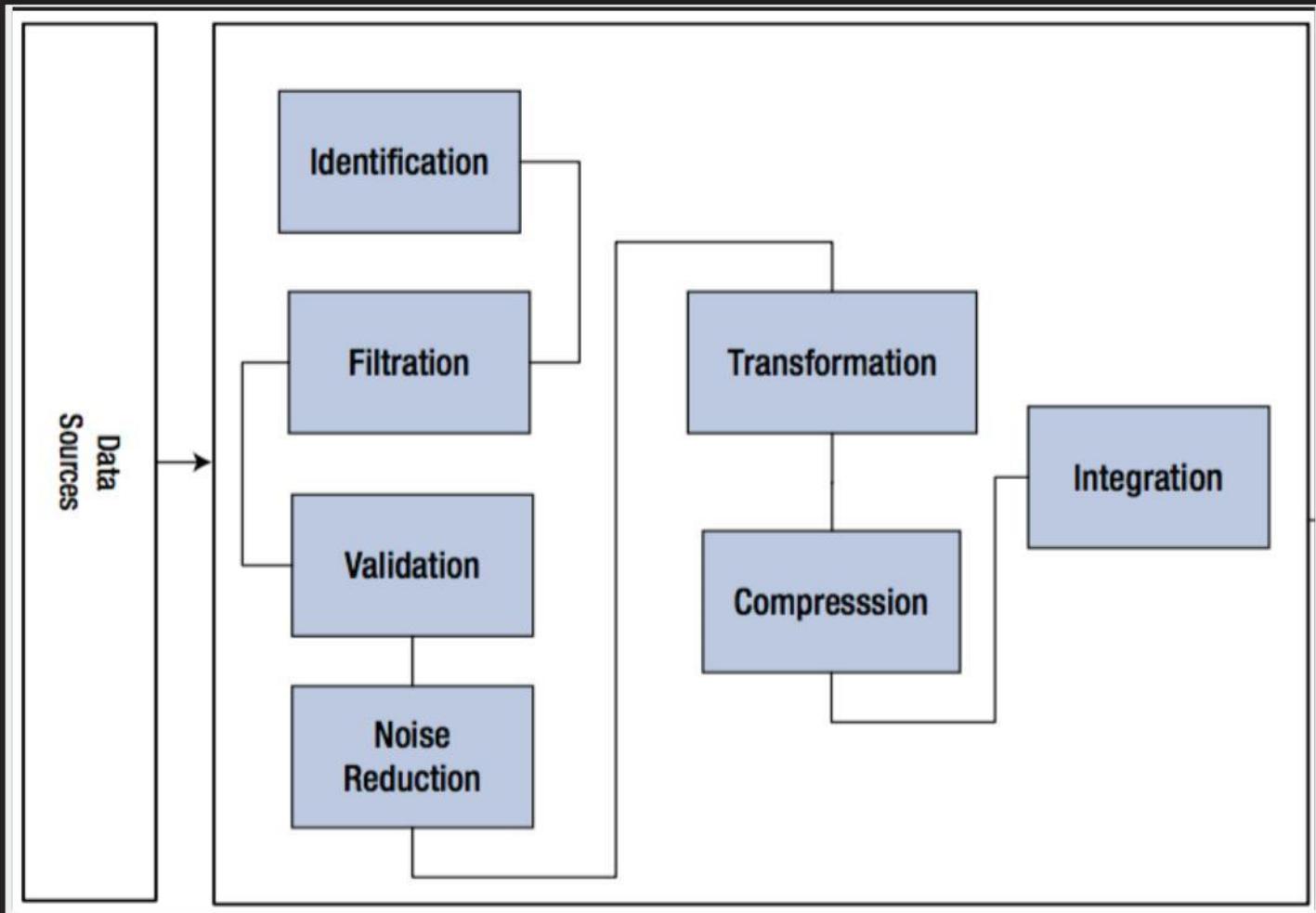
# Ingestion Layer

**Signal-to-Noise  
ratio  
10:90**

separate the  
noise from  
relevant info

# It has capability to

- Validate
- Cleanse
- Transform
- Reduce
- Integrate



# Distributed Storage Layer

# Fault tolerance Parallelization

# HDFS

## massively scalable distributed file system

# HDFS

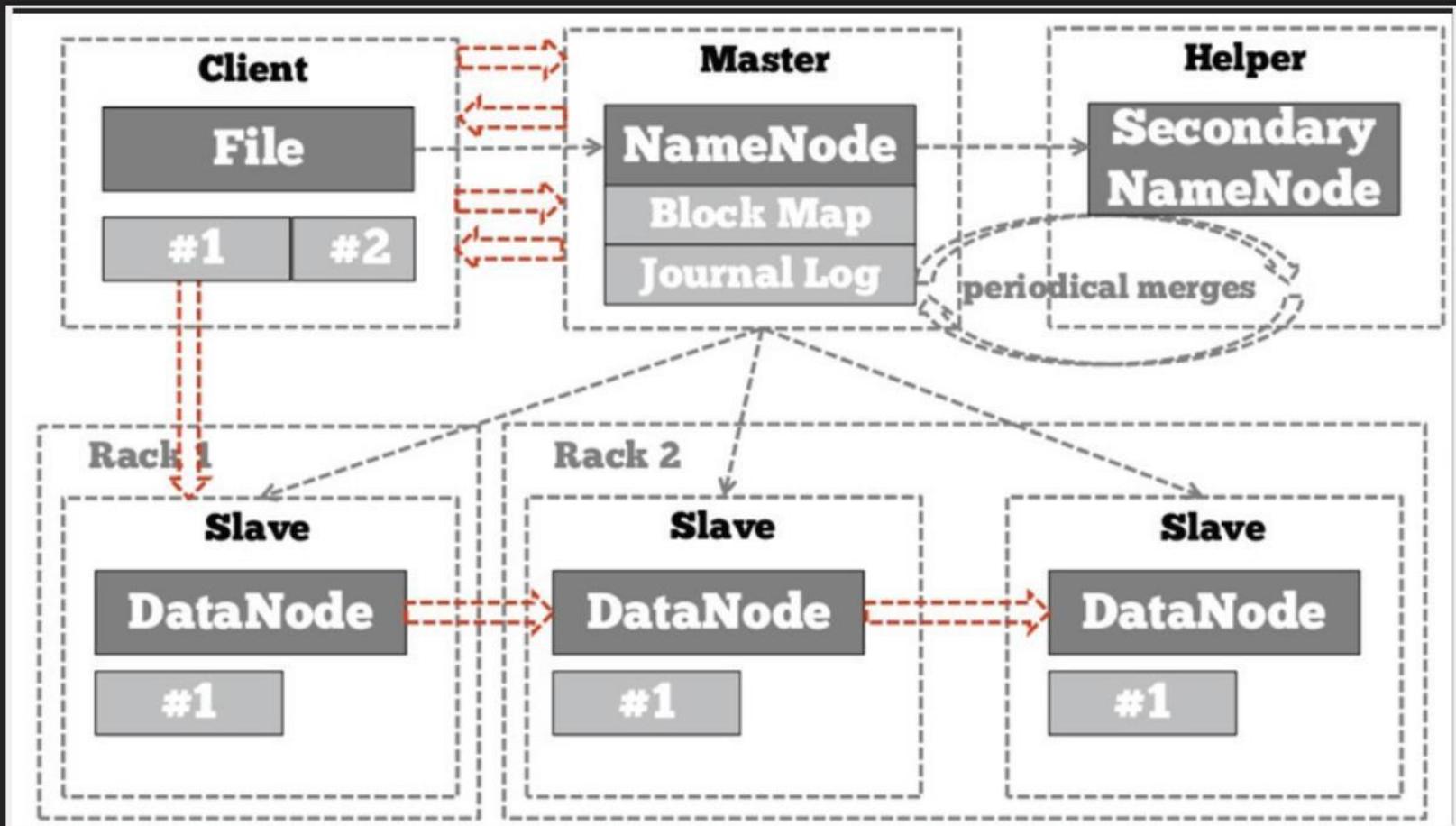
**Intended for**

- **large files**
- **batch inserts**

*Write once, read many times*



# HDFS Architecture

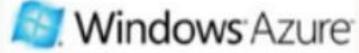
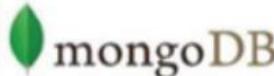


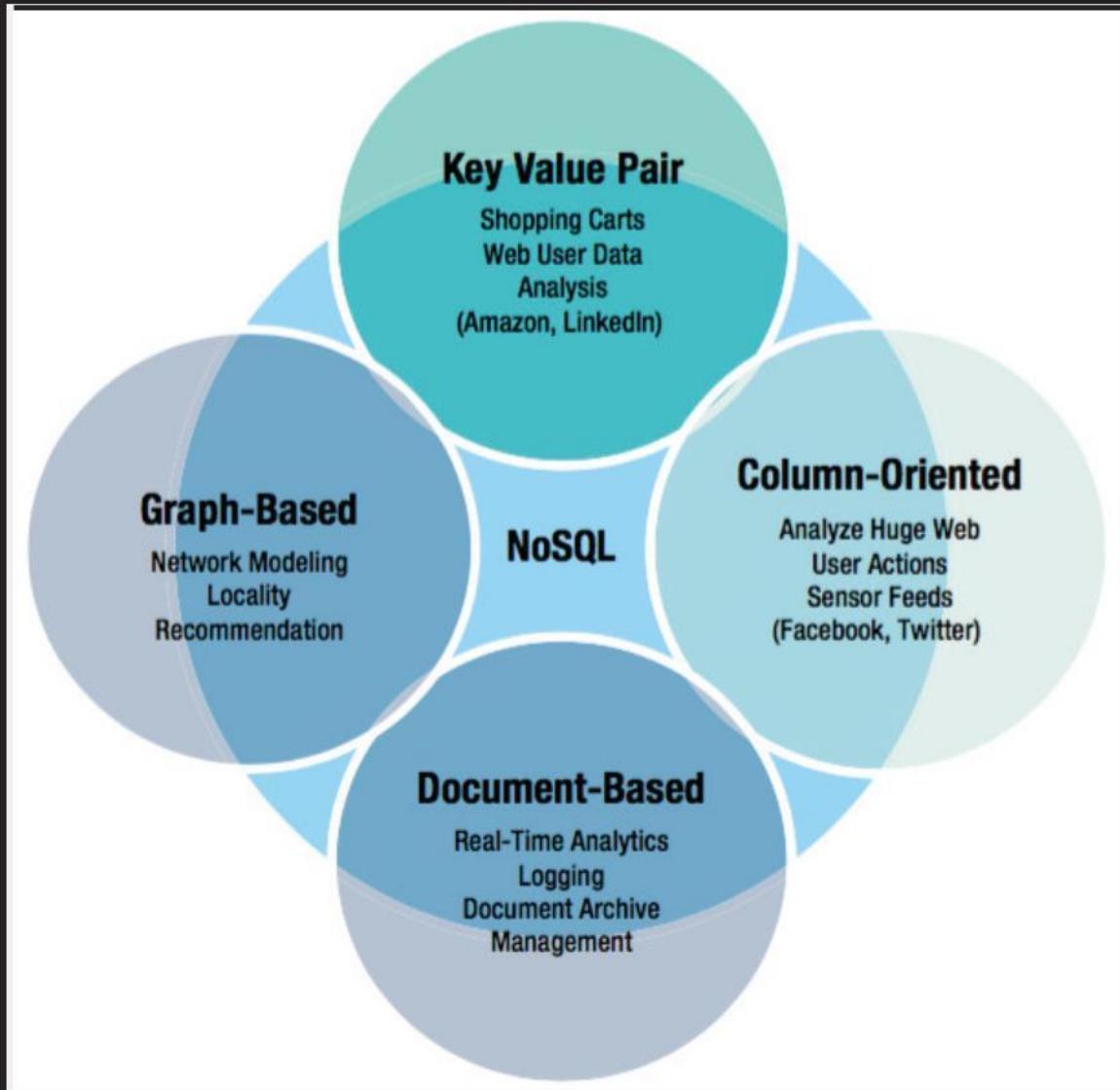
Non-relational,  
distributed data?

# NoSQL

# CAP theorem

## Consistency, Availability, Partition Tolerance

Key-Value Data Stores	Column-oriented Data Stores	Document Data Stores	Graph Data Stores
    	     	  	    



# Ingestion to DFS

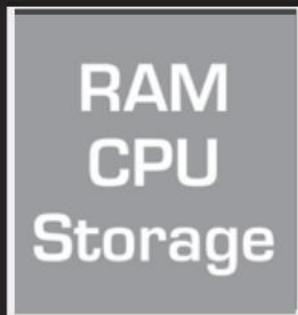
## Sqoop, Flume, MapReduce, ETL

# Infrastructure & Platform Layer

# Computing & Scalability

# Hadoop?

# Vertical Scaling



# Vertical Scaling

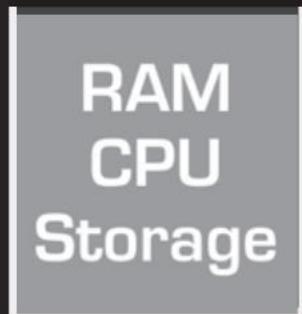


RAM  
CPU  
Storage

# Vertical Scaling

RAM  
CPU  
Storage

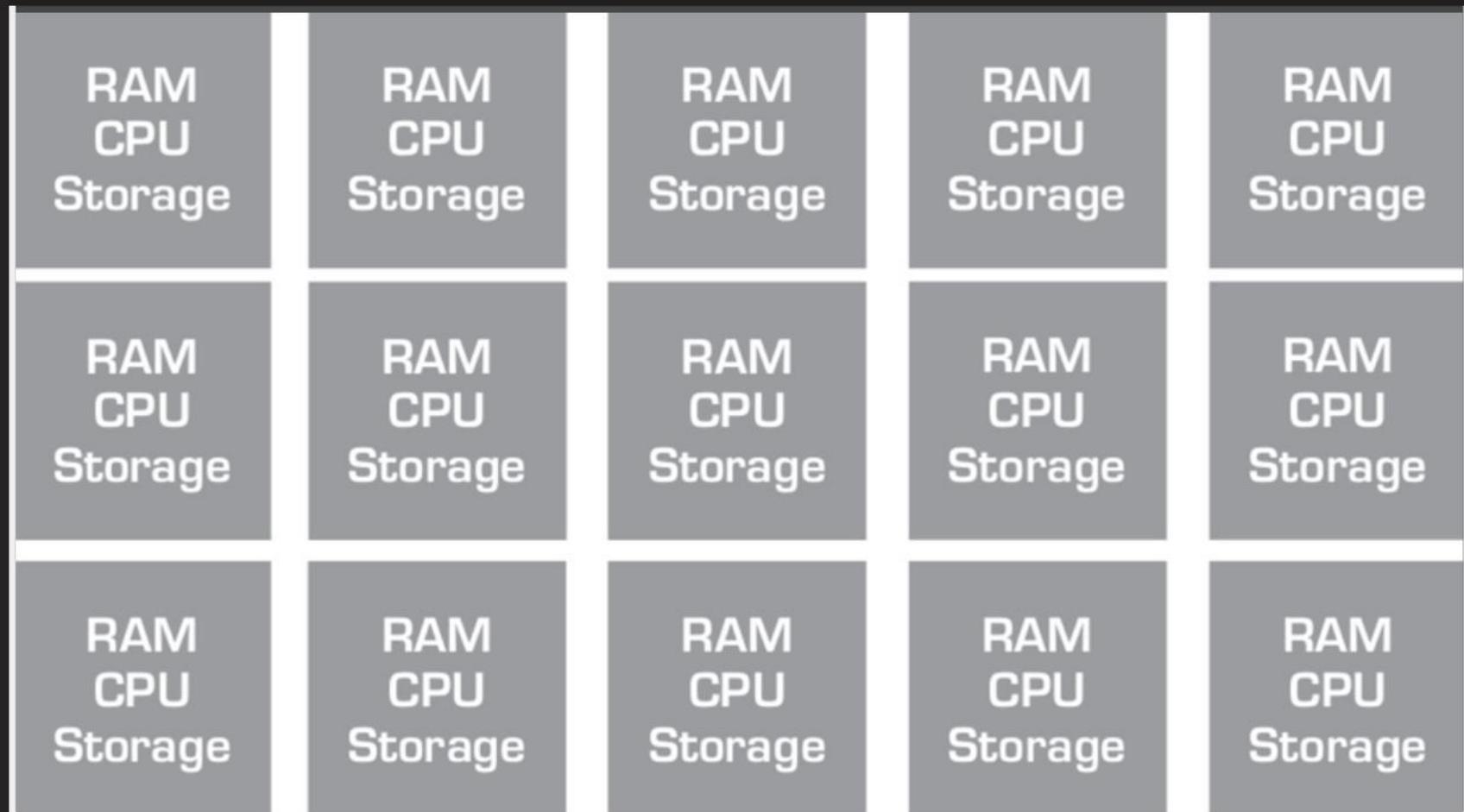
# Horizontal Scaling

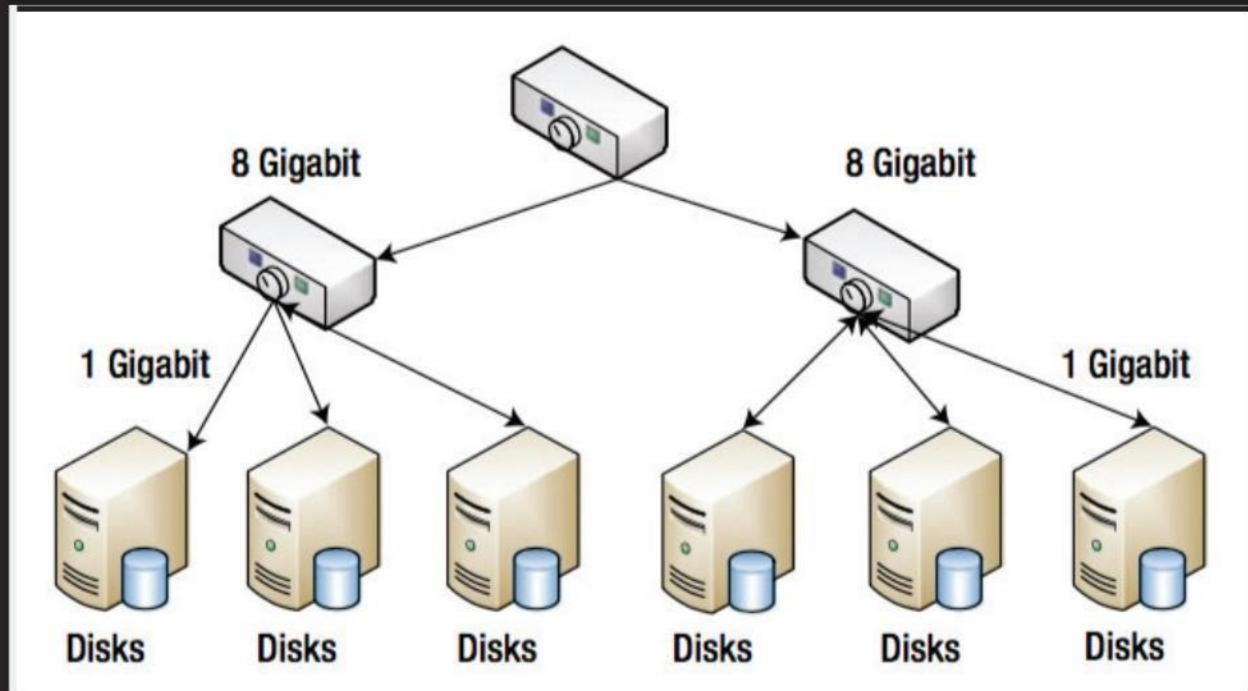


# Horizontal Scaling



# Horizontal Scaling

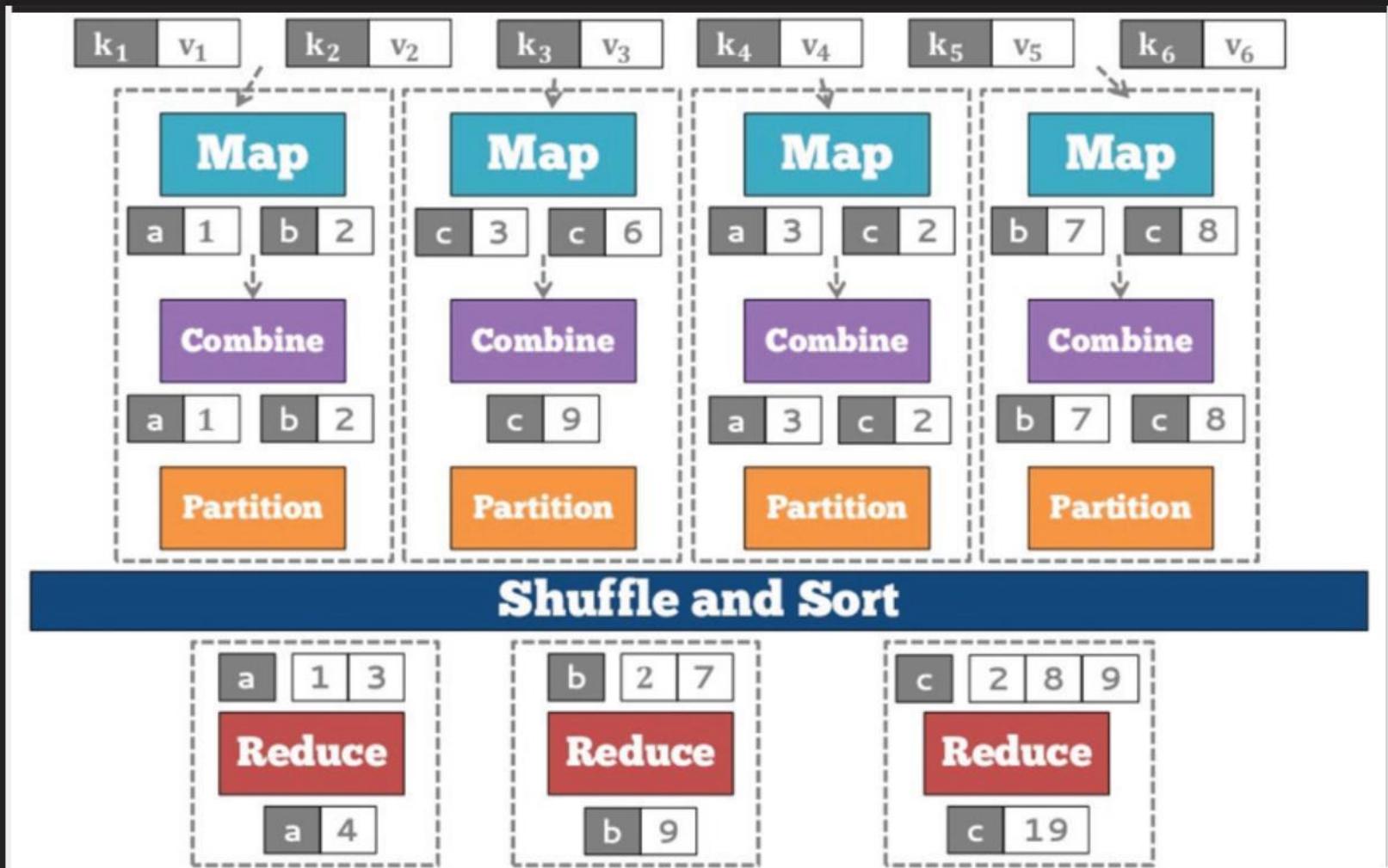


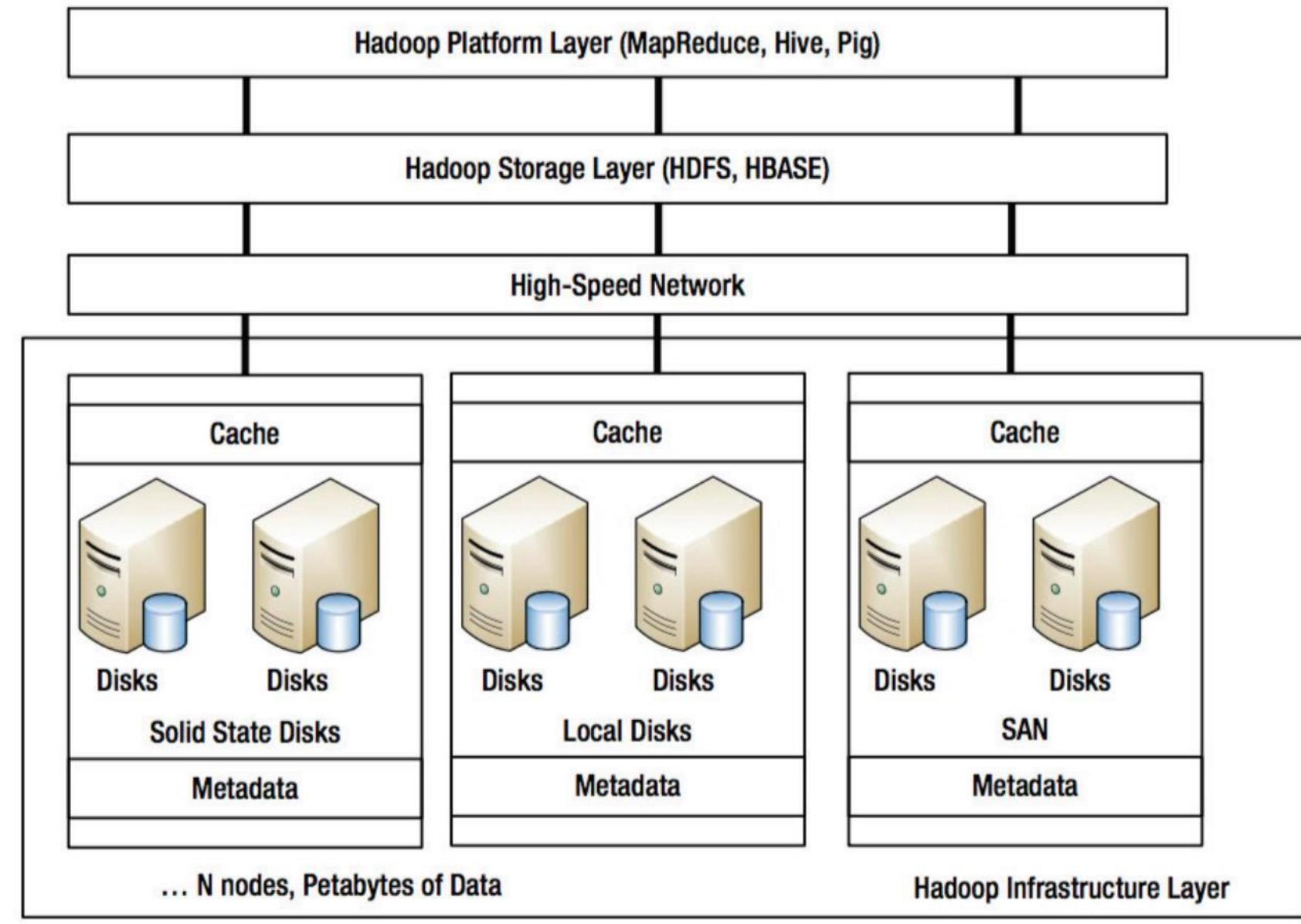


# MapReduce

is the main computation paradigm

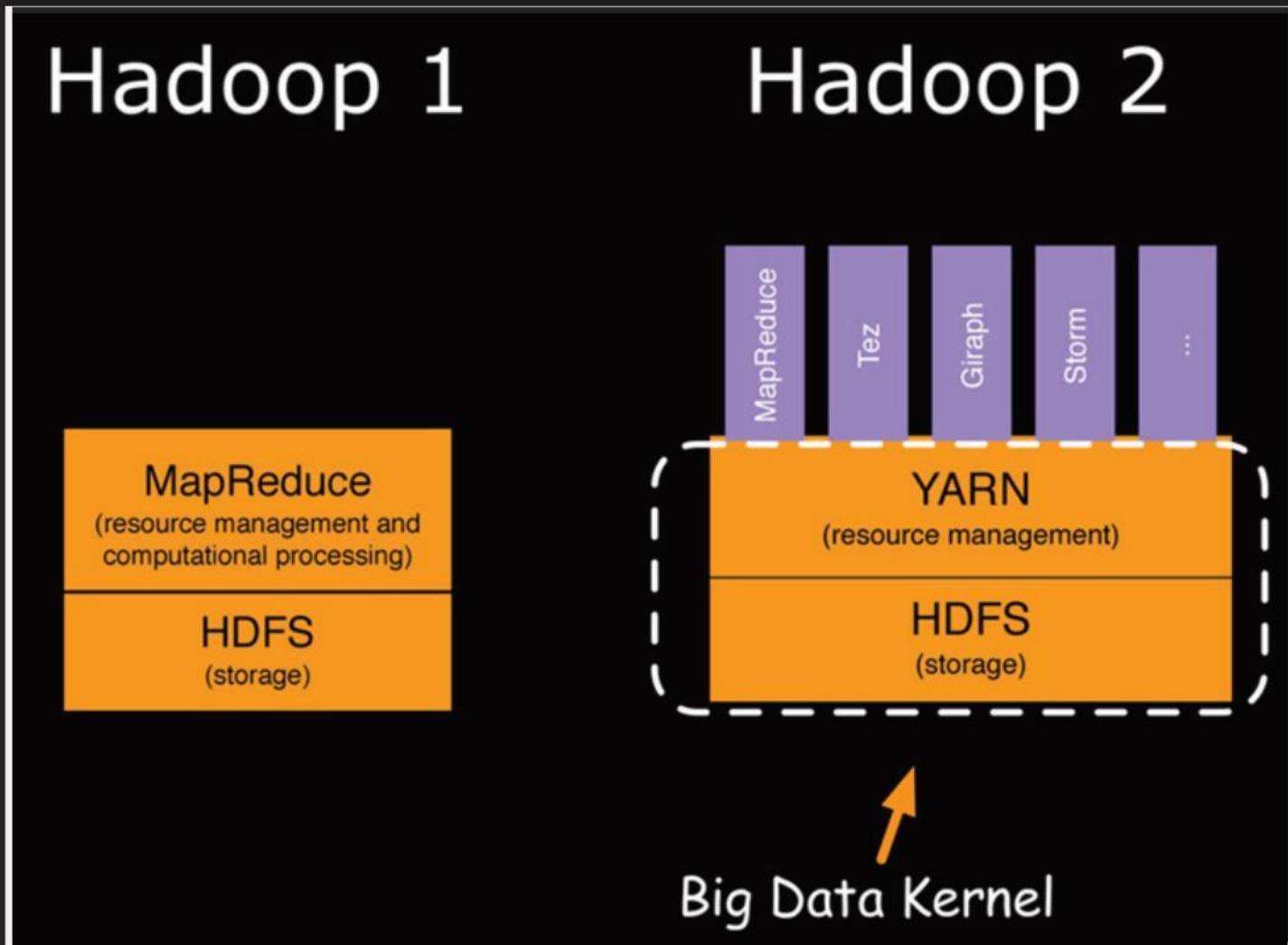
# MapReduce



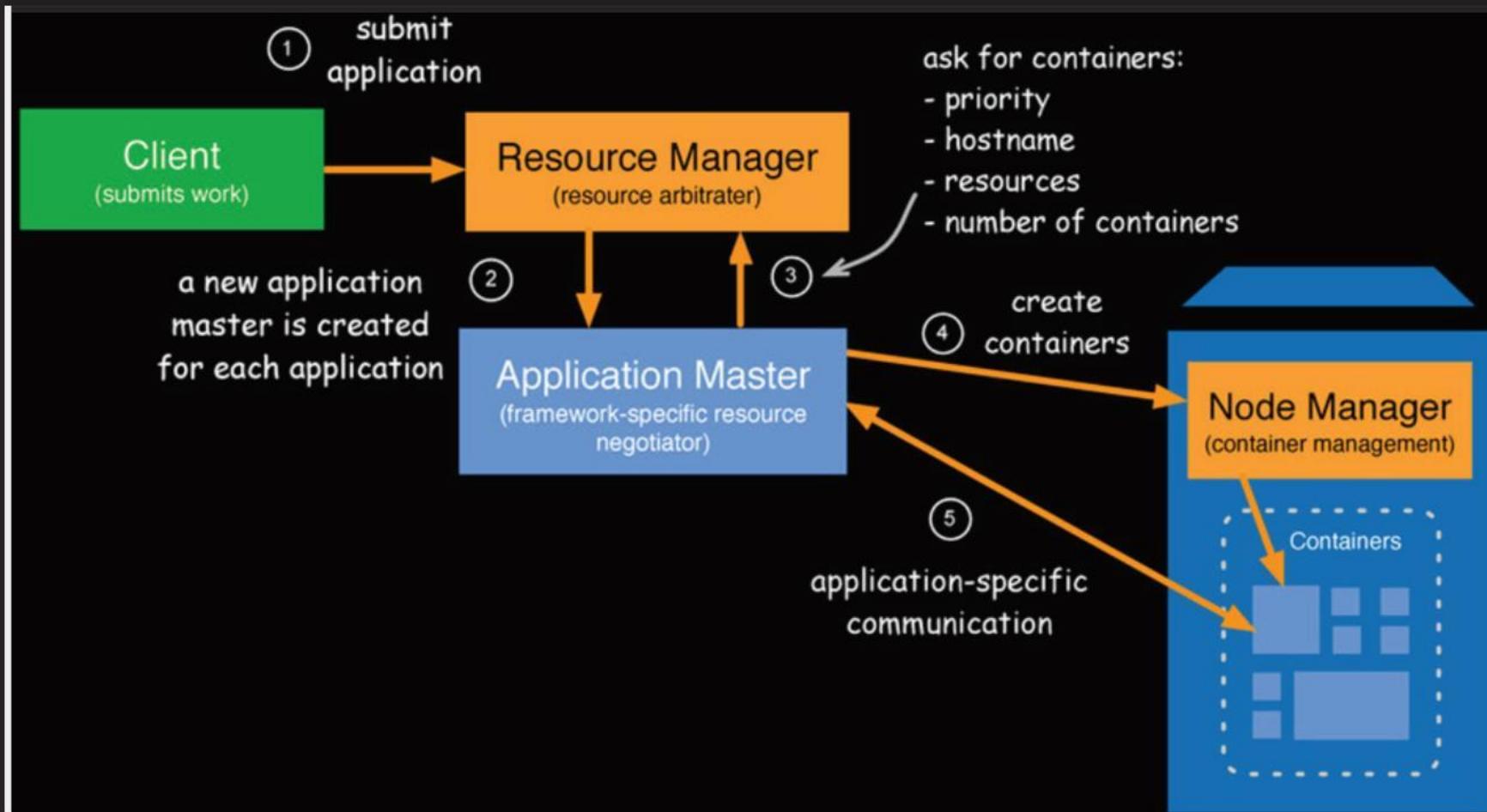


# Hadoop 2

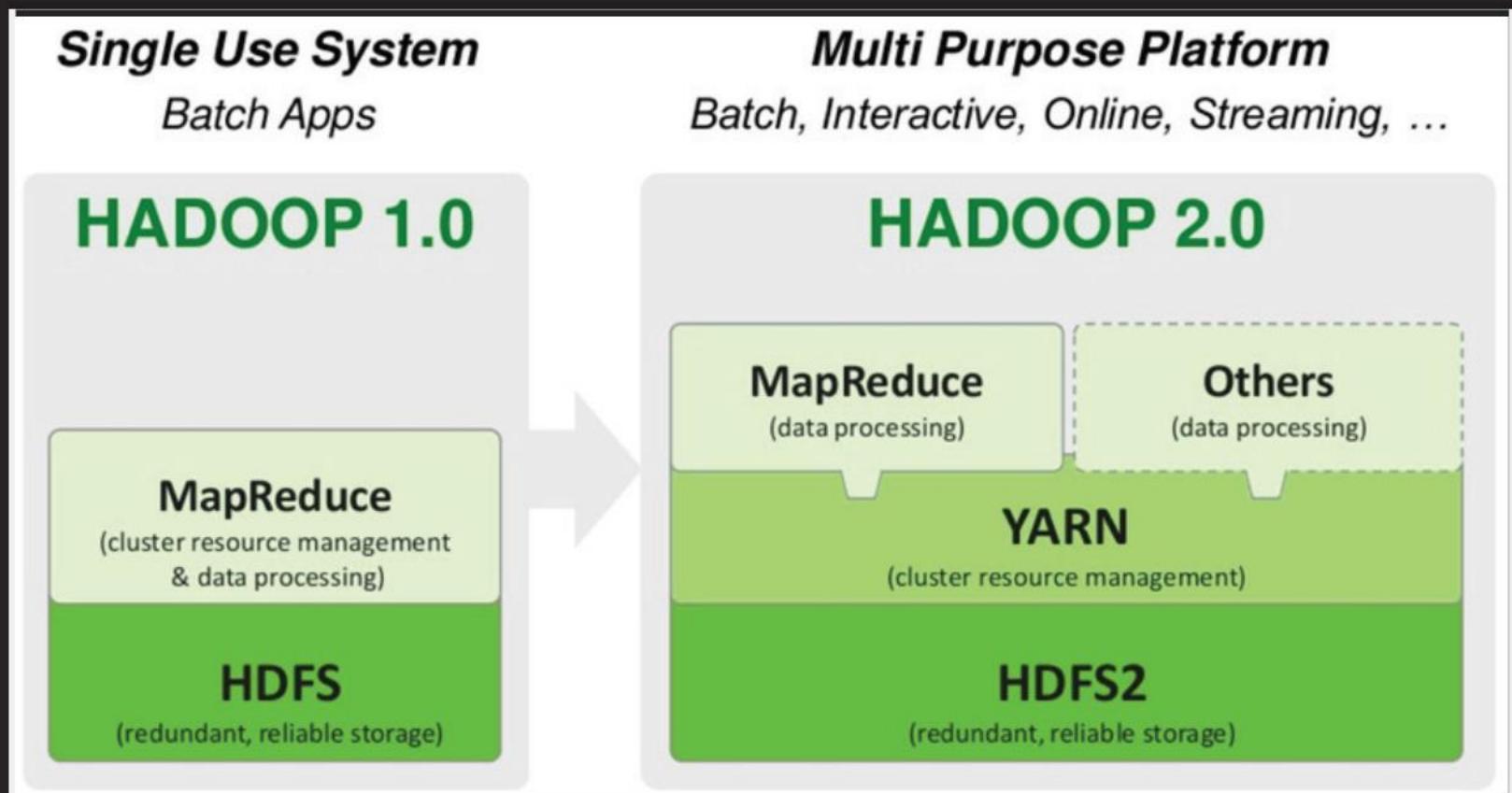
# What's new?



# What's new?



# H1 vs. H2



**One cluster,  
distributed storage,  
distributed scheduler,  
many types of applications.**

# Blueprints

- NoSQL with HBase
- Stream Processing with Storm/Spark
- Graph Processing with Giraph
- SQL on Hadoop with Impala
- Columnar Data Formats

# Security Layer

# Data need to be **protected**

- Meet compliance requirements
- Individual's privacy

Proper  
authorization and  
authentication  
needed

# What can we do?

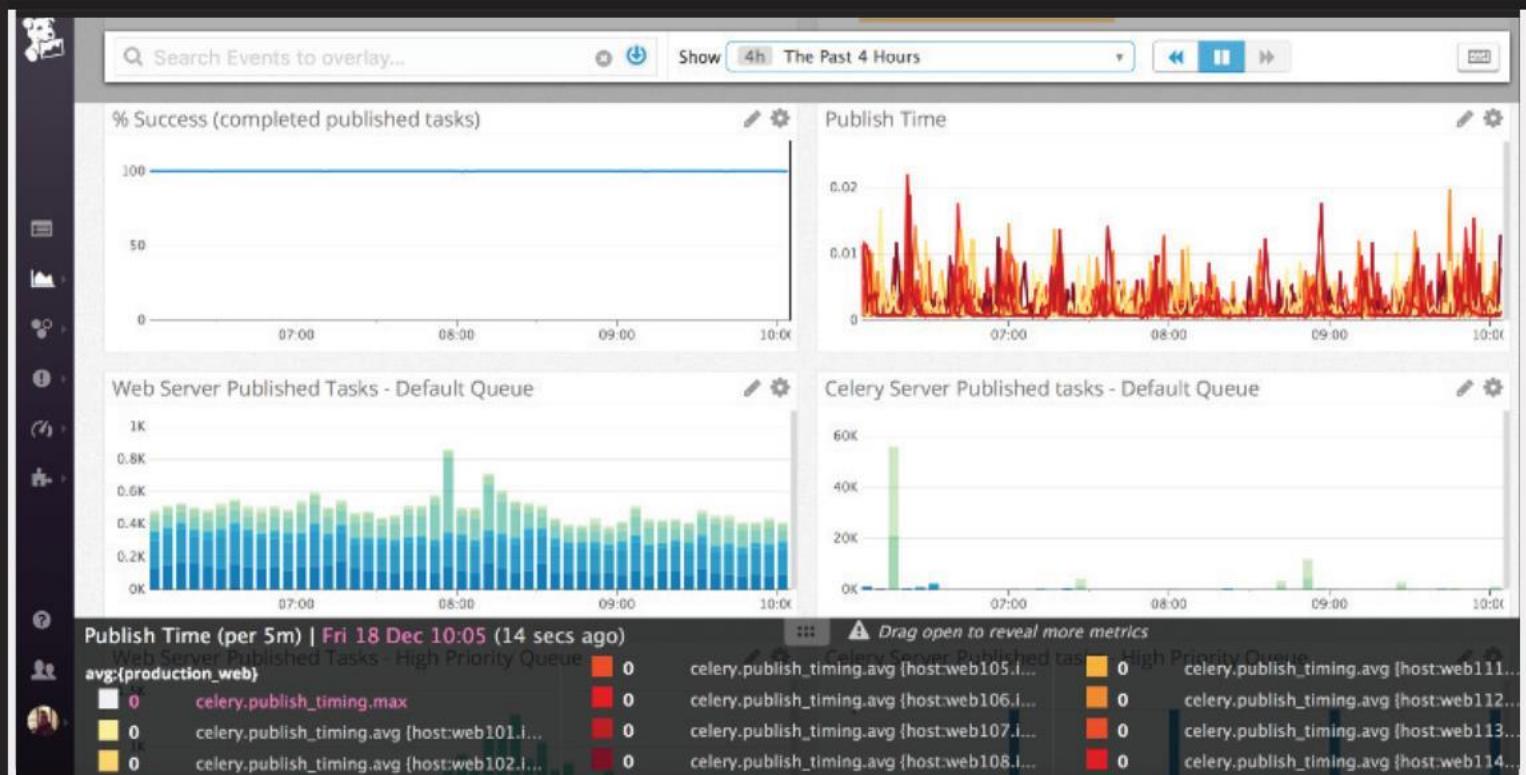
- Authentication protocol like Kerberos
- Enable file layer encryption
- Use SSL, certificates and trusted keys
- Provision with Chef, Puppet or Ansible like tools
- Log all the communication for detecting anomalies
- Monitor whole system

# Monitoring Layer

**Get a complete  
picture  
of our Big Data tech stack**

Satisfy SLAs with  
min downtime

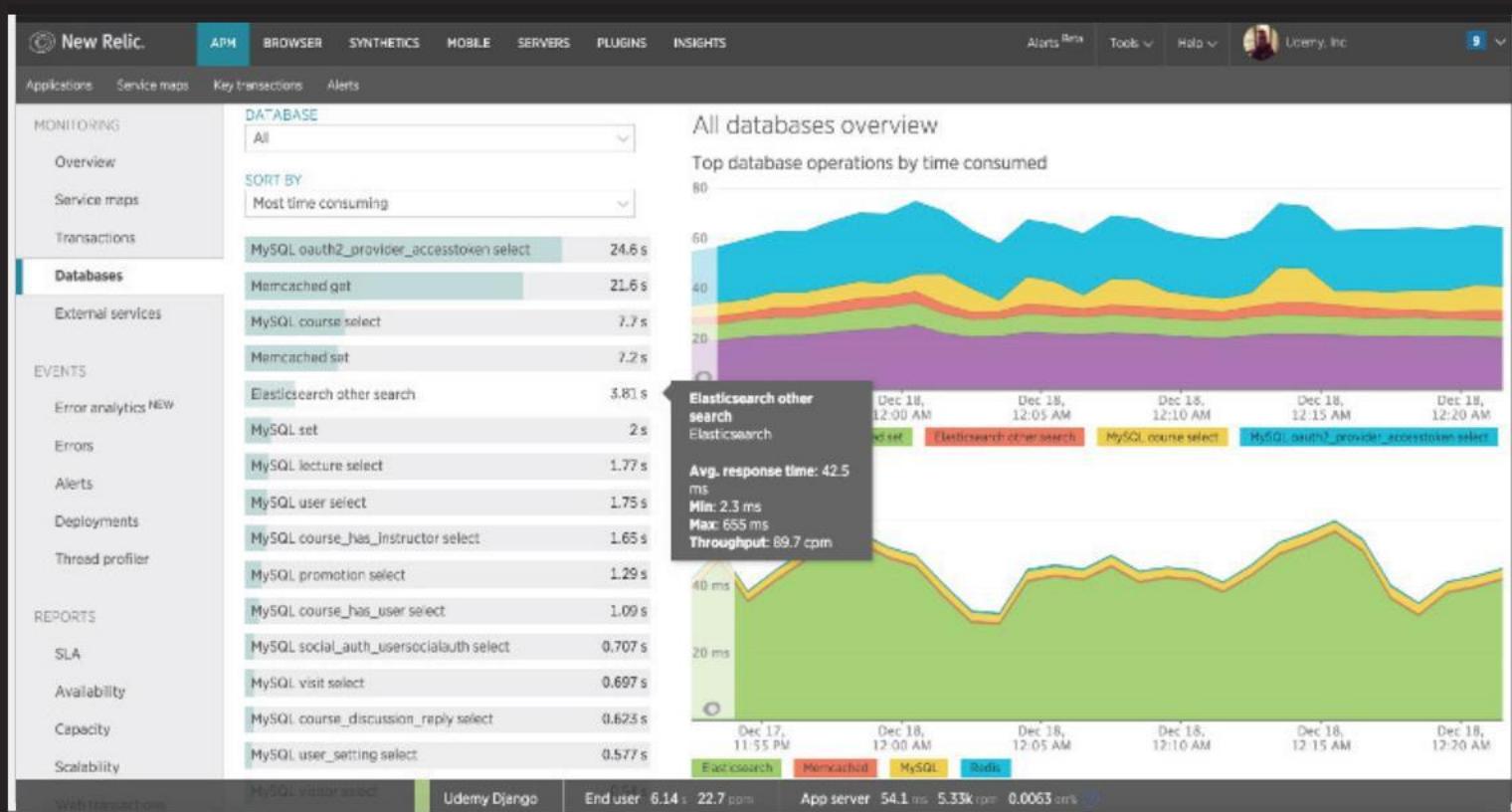
# DataDog



# New Relic (Overview)



# New Relic (Databases)



# Analytics Engine

# Co-Existence with Traditional BI

- Data warehouse in the traditional way
- Distributed MR processing on big data stores

**Mediate data in either direction**  
**i.e use **Hive/HBase** with Sqoop**

**Real-time analysis can leverage  
low-latency NoSQL stores**

i.e Cassandra, Vertica, ...

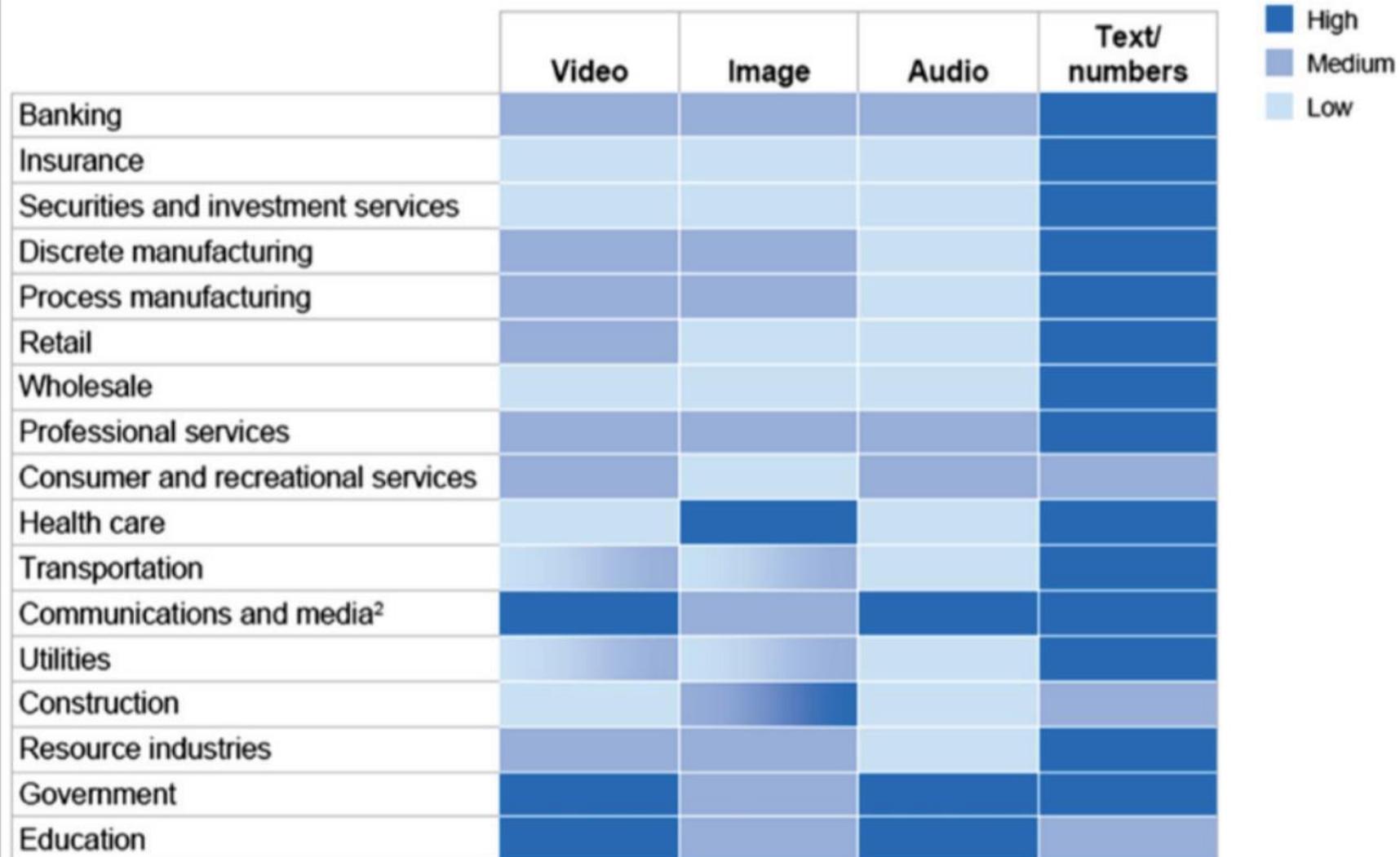
**R may be used for complex  
statistical algorithms**

# Search Engines

Huge volume and  
variety of data

**“needle in a  
haystack”**

## The type of data generated and stored varies by sector<sup>1</sup>



1 We compiled this heat map using units of data (in files or minutes of video) rather than bytes.

2 Video and audio are high in some subsectors.

Need blazing **fast** search  
mechanism  
to **index** and **search** for big data  
analytics

## **Result Display**

## **Query Processing**

## **User Management**

## **Search Functions**

Spelling

Stemming

Faceting

Highlighting

Tagging

Parsing

Semantics

Pertinence

## **Search Engine**

Indexing

Crawling

Elastic Search,  
Solr, ...

# Real-time Processing

# In memory?

**Apache Spark™** is a fast and general engine for large-scale data processing.

Apache Spark

Spark  
SQL

Spark  
Streaming

MLlib  
(machine  
learning)

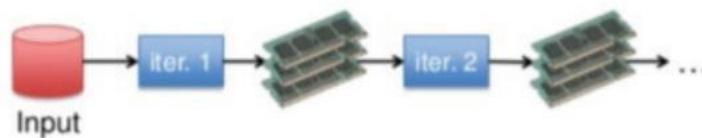
GraphX  
(graph)

Apache Spark

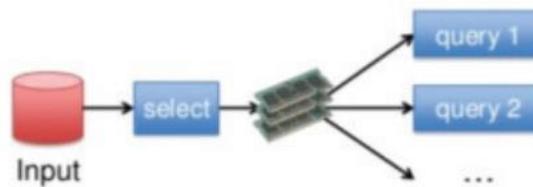


## Spark In-memory Processing

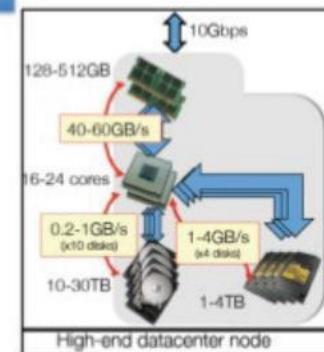
Iterative:



Interactive:



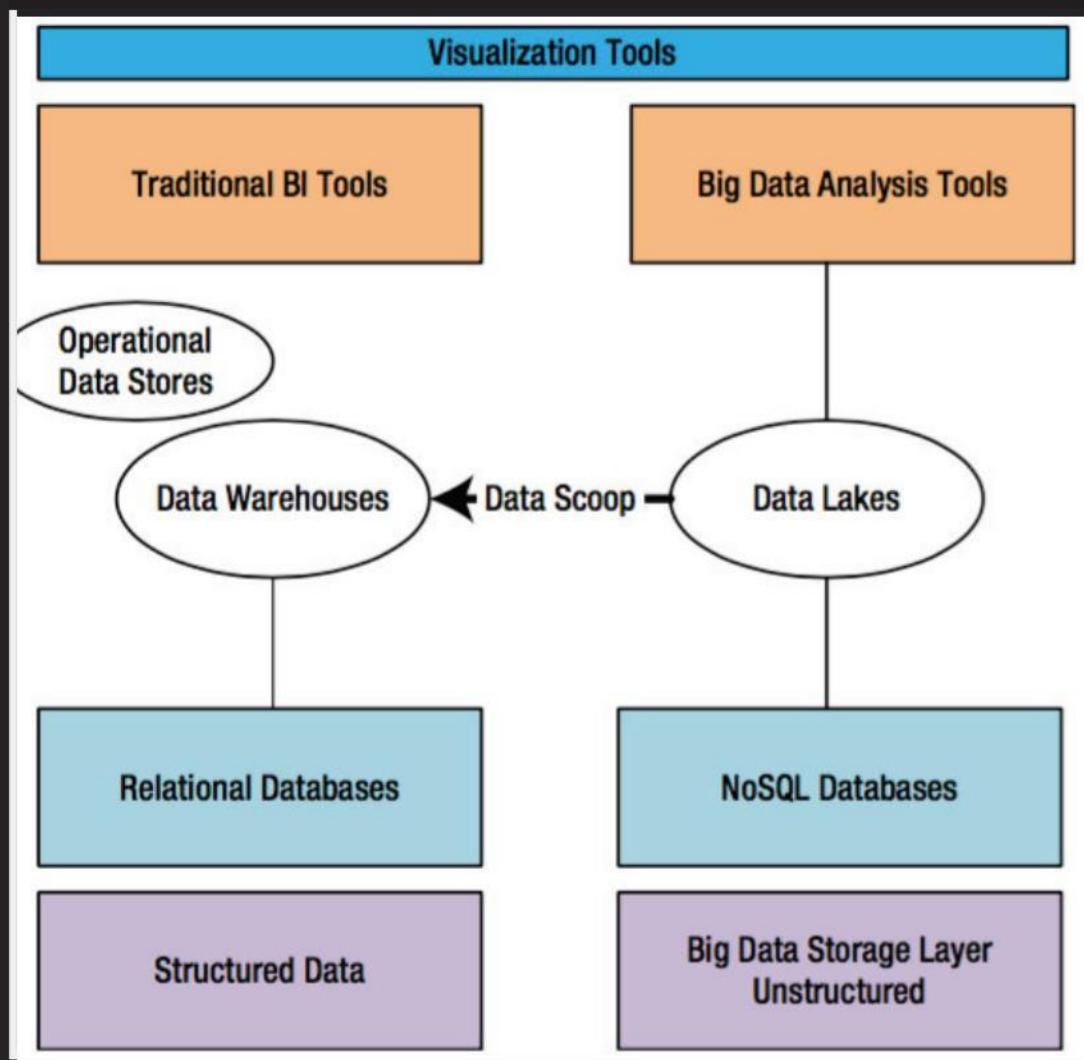
- 1.Extract a working set
- 2.Cache it
- 3.Query it repeatedly



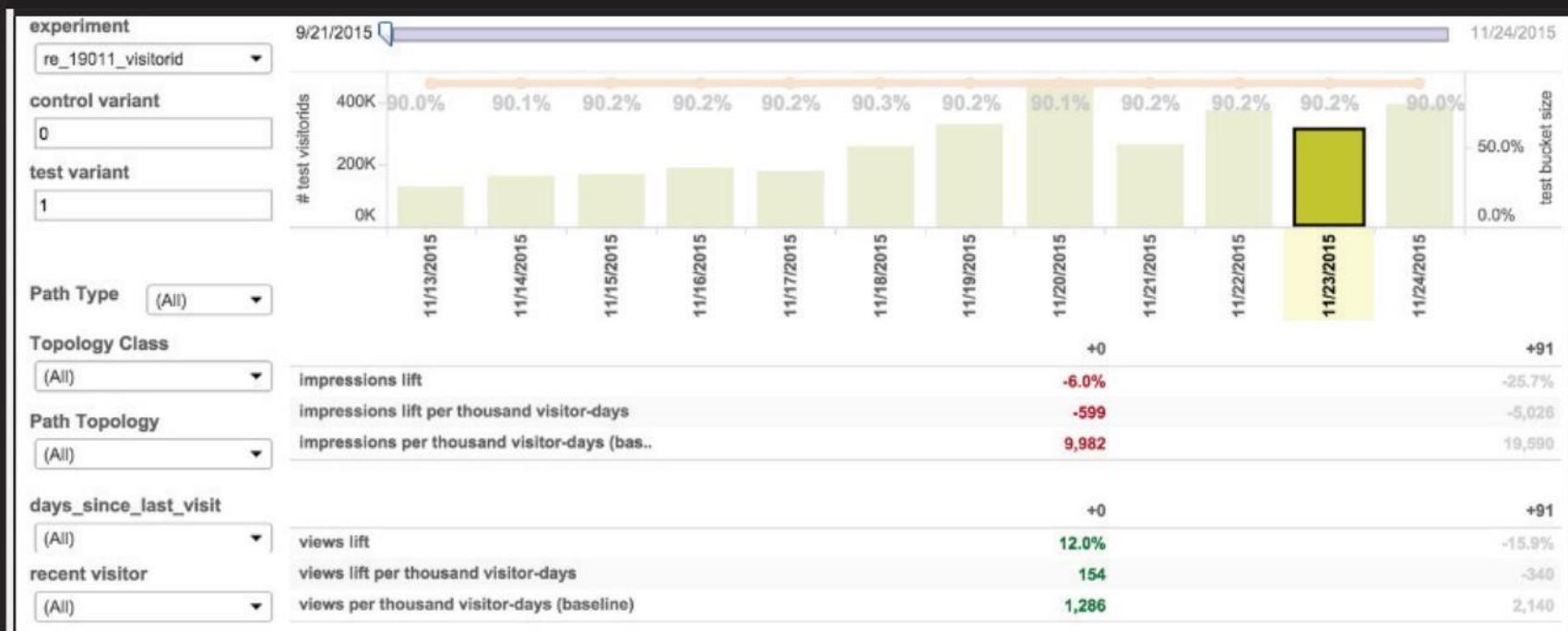
**Storm, Kinesis,  
Flink, ...**

# Visualization Layer

**Gain insight faster**  
**Look at different aspects of**  
**data visually**



# Tableau



# ChartIO

## Learning at Udemy

**357,772**

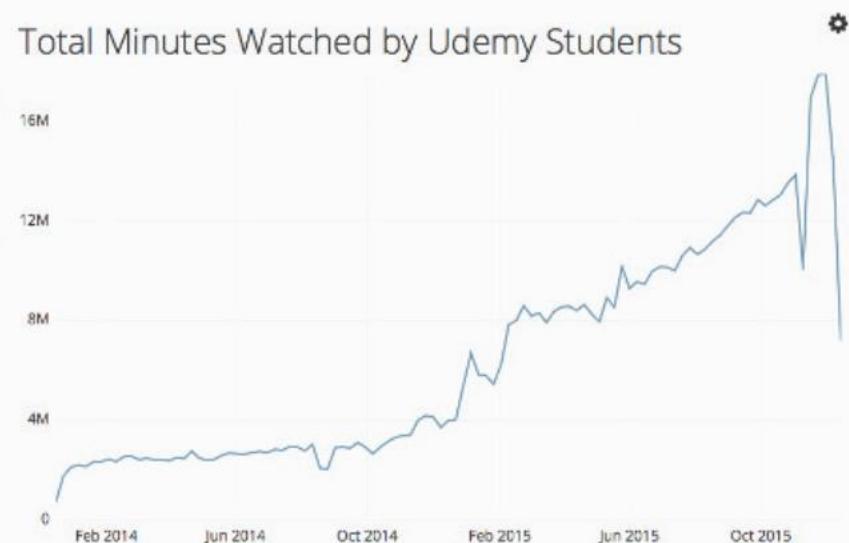
Days of Content Consumed  
in 2015

*Students consumed an incredible 352,000 days of content on Udemy in 2015; average consumption also went up across the platform.*

monthly average min consumed

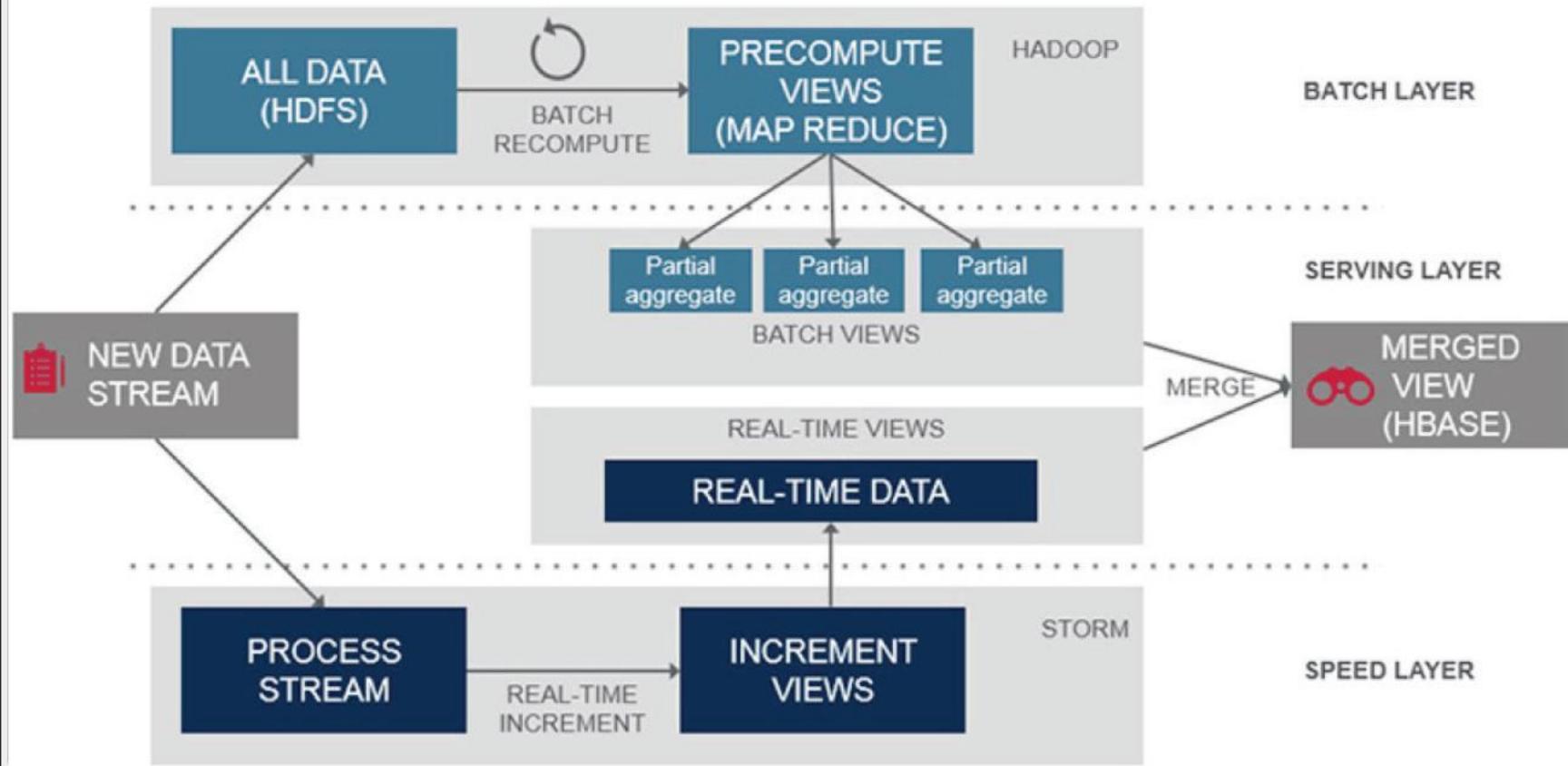


Total Minutes Watched by Udemy Students



# Lambda Architecture

# Lambda Architecture



Lambda Architecture / MapR

# Don't forget

There is no  
"One Size Fits All"  
solution

We need  
**Continuous  
Development**



*Get Big*  
by starting

*small*

A close-up photograph of a person's hands holding a black Nikon DX NIKKOR 18-55mm f/3.5-5.6G VR lens. The lens barrel is dark with white text. In the center, it says "NIKKOR". Around the center, it says "18-55mm f/3.5-5.6G VR". At the bottom, it says "DX NIKKOR". The background is blurred.

Focus on  
*Business Impact.*

# Thank You :)