# What Is Big Data?

# Volume of data created, captured, copied, and consumed worldwide

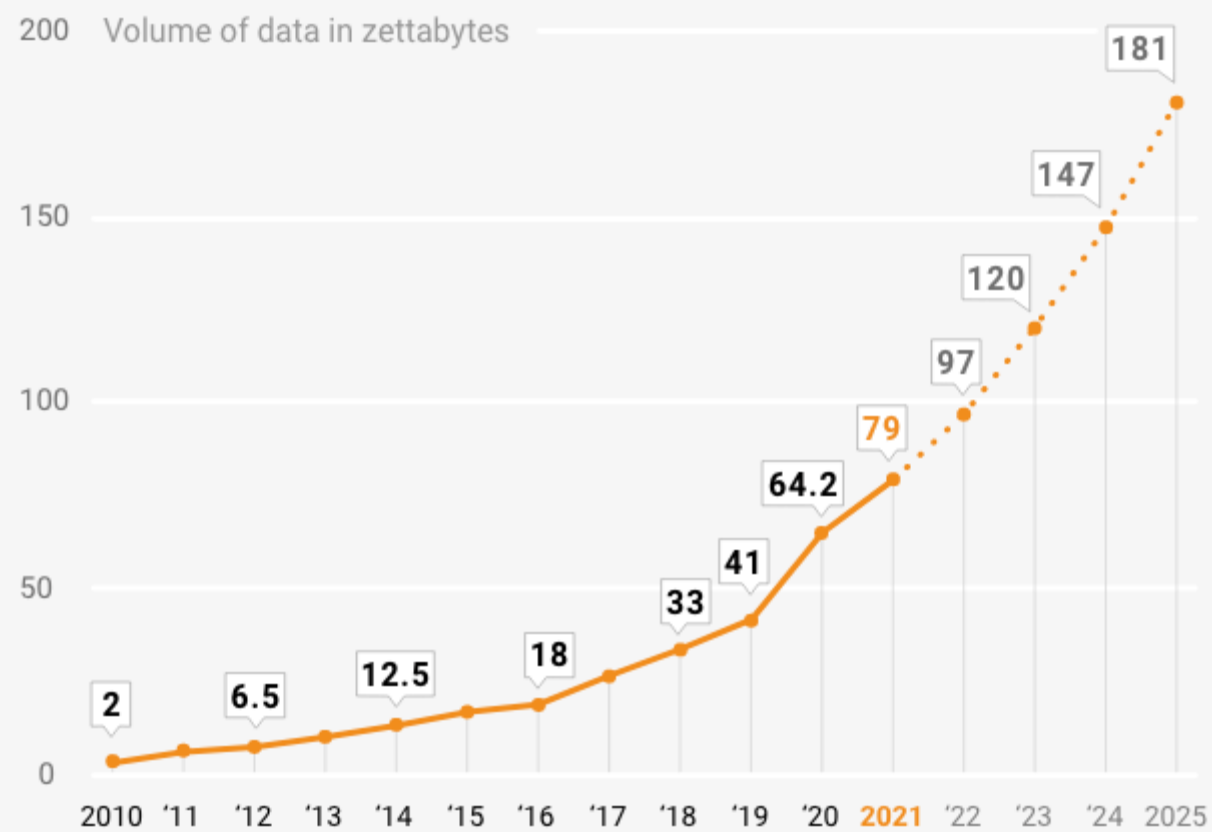The volume of data generated, consumed, copied, and stored is projected to exceed 180 zettabytes by 2025

**Volume of data in zettabytes**

200

181

147

150

120

100

97

79

64.2

50

41

33

18

12.5

6.5

2

2010 '11 '12 '13 '14 '15 '16 '17 '18 '19 '20 2021 '22 '23 '24 2025

Source: statista.com

---

**1** How much data is generated every minute?

Source: Domo

**41,666,667**
messages shared
by WhatsApp users

**1,388,889**
video / voice calls made
by people worldwide

**404,444**
hours of video streamed
by Netflix users

**347,222**
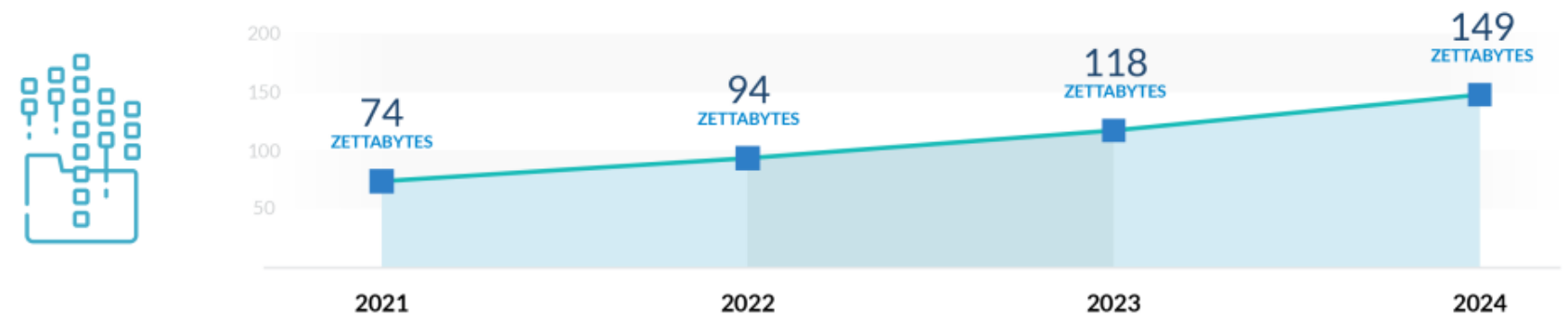stories posted by Instagram users

**150,000**
messages shared by Facebook users

**147,000**
photos shared by Facebook users

**2** Estimated Data Consumption from 2021 to 2024

Source: IDC / Statista

200

149
ZETTABYTES

150

118
ZETTABYTES

94
ZETTABYTES

100

74
ZETTABYTES

50

2021    2022    2023    2024

**3** Data Growth in 2021

Sources: TechJury, Internet Live Stats, Cisco, PurpleSec

**2 TRILLION**
searches on Google by the end of 2021

**1.134 TRILLION MB**
volume of data created every day

**3,026,626**
emails sent every second, 67% of which are spam

**278,108 PETABYTES**
global IP data per month by the end of 2021

**230,000**
new malware versions created every day

**82%**
share of video in total global internet traffic at the end of 2021

# IS THERE REALLY A USE CASE?

### Science

- Large Hadron Collider - 1 Petabyte every second
- NASA - 1.73 Gigabyte every hour

### Government

- NSA - Utah Data Center - Yottabyte Capacity
- Big Data Research and Development Initiative
- Barack Obama's successful 2012 re-election campaign

### Private

- eBay - 40PB Hadoop cluster for search, consumer recommendations, and merchandising
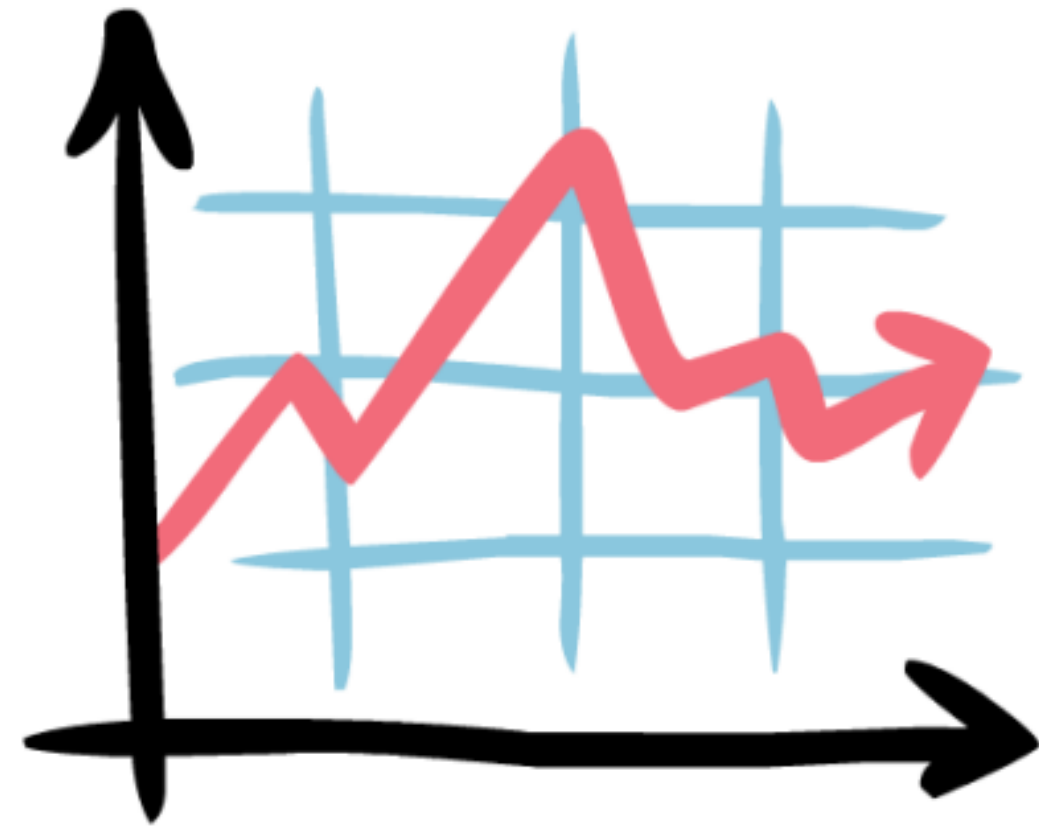- Facebook - 30 PB Hadoop cluster. 50 billion photos. 130TB of logs every day.

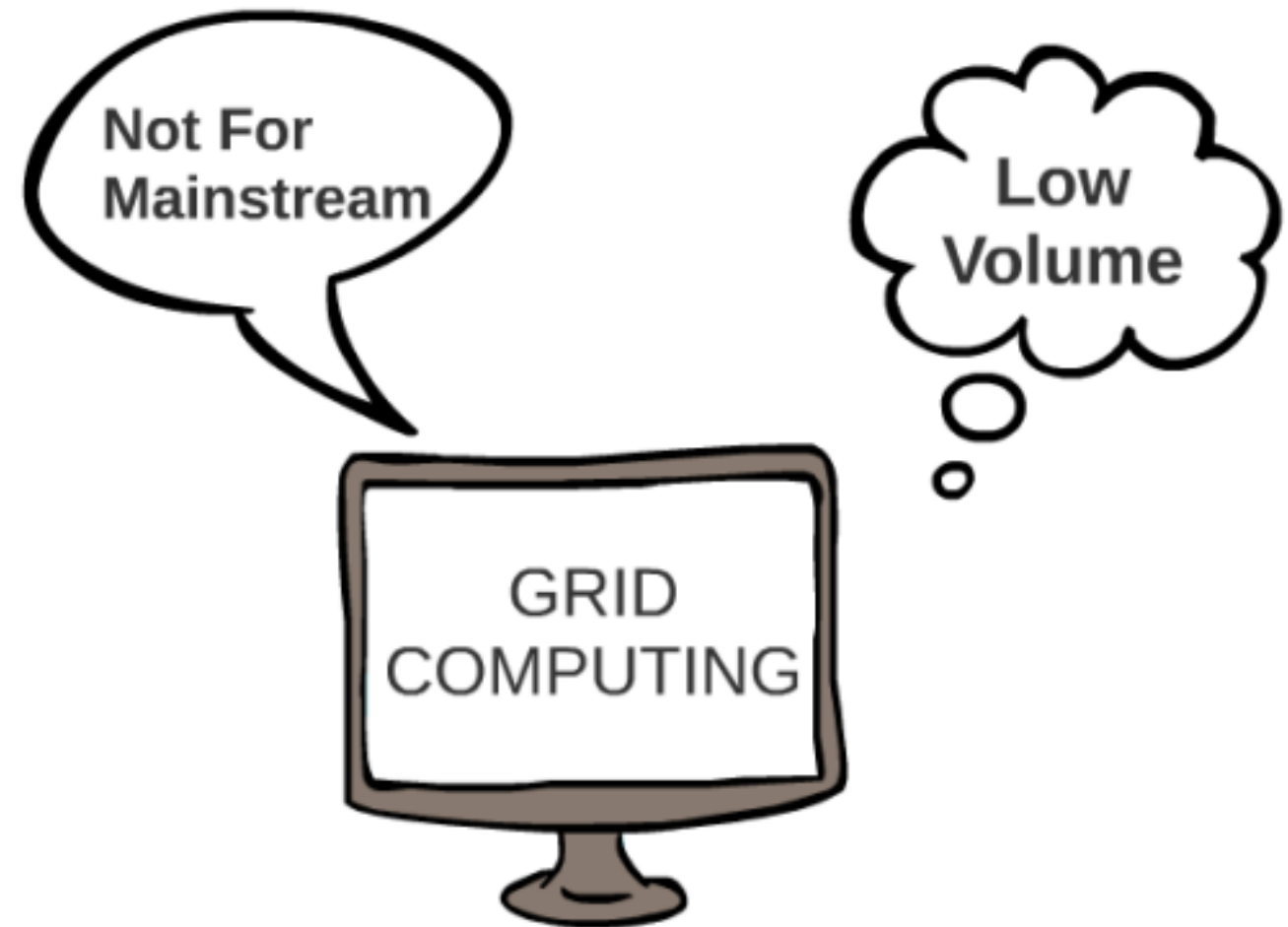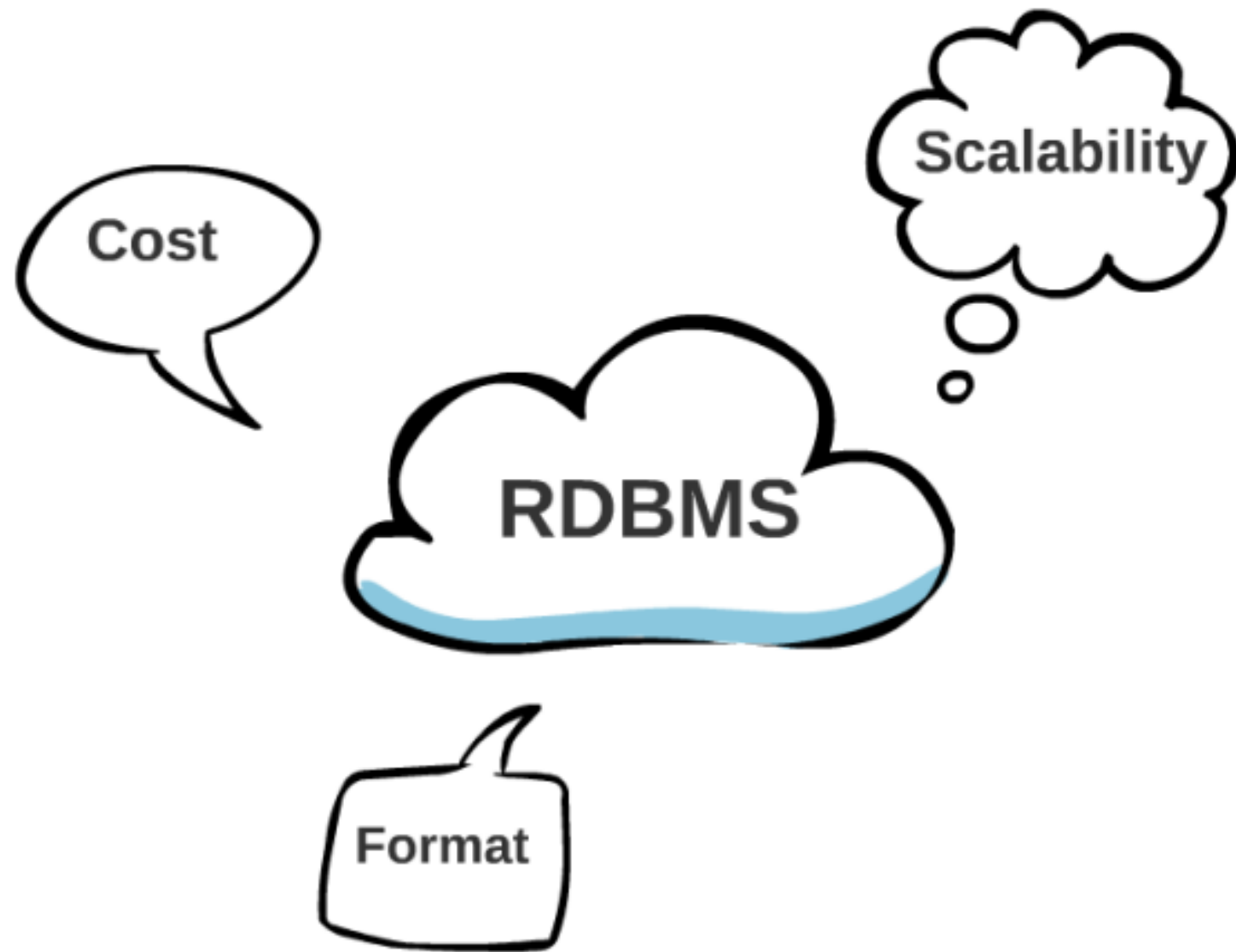# BIG DATA - CHALLENGES

Storage

Computational Efficiency

Data Loss

Cost

# HADOOP - A GOOD SOLUTION

✓ Support Huge Volume

✓ Storage Efficiency

✓ Good Data Recovery Solution

✓ Horizontal Scaling

✓ Cost Effective

✓ Easy For Programmers & Non Programmers
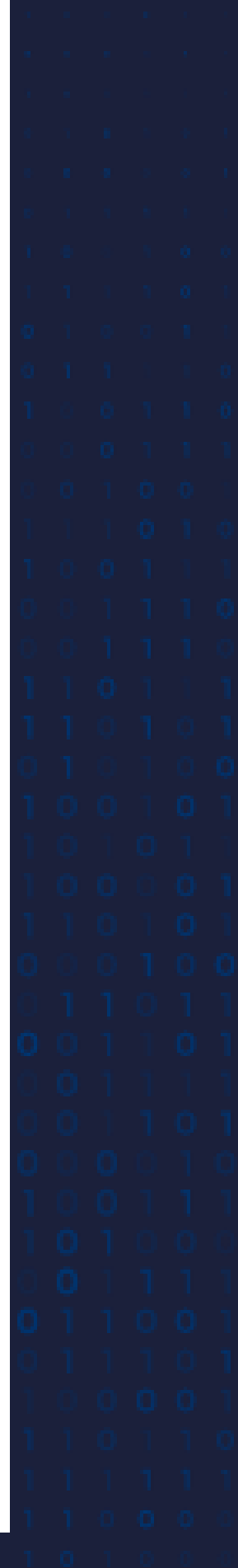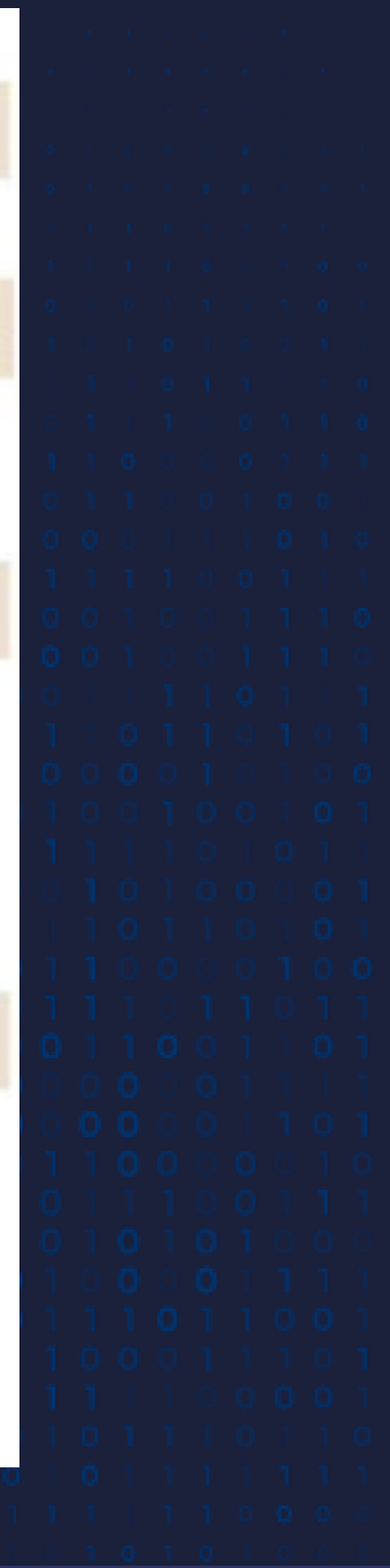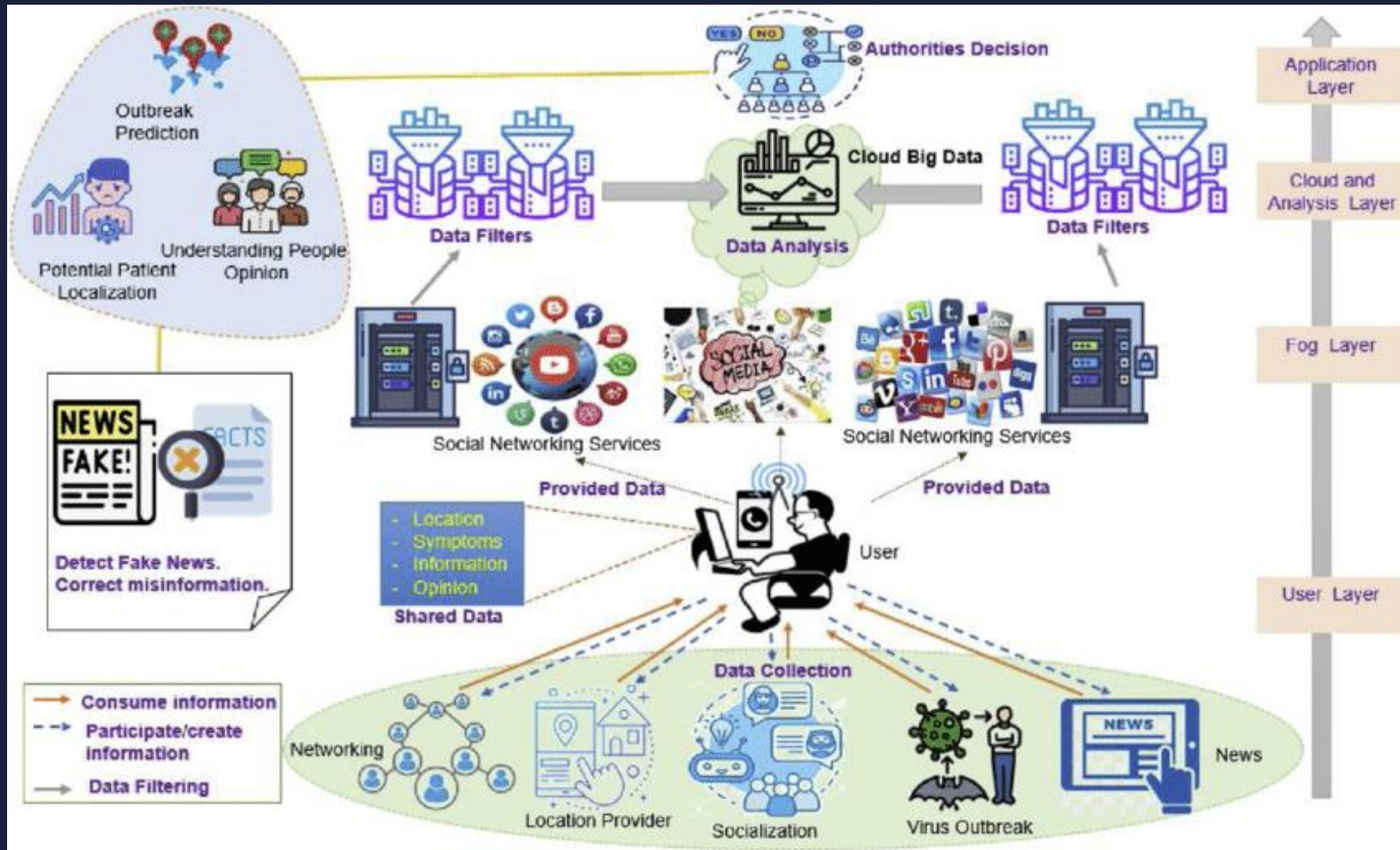
# Big Data Application

**Business Intelligence and Analytics:**

Big data is used to analyze historical and real-time data to identify trends, patterns, and correlations, helping organizations make informed decisions, optimize operations, and develop data-driven strategies.

**Customer Insights:**

Analyzing vast amounts of customer data, including social media interactions, purchase history, and demographic information, helps businesses understand customer behavior and preferences, enabling targeted marketing and improved customer experiences.

**Fraud Detection and Security:**

Big data analytics can be employed to detect fraudulent activities and enhance cybersecurity by identifying anomalies and patterns indicative of cyber threats.

## Healthcare Analytics:

Analyzing electronic health records, medical imaging data, and genomic information can lead to improved patient care, disease prediction, and drug discovery.

## Predictive Maintenance:

In industries like manufacturing and aviation, big data is used to predict equipment failures and optimize maintenance schedules, reducing downtime and costs.
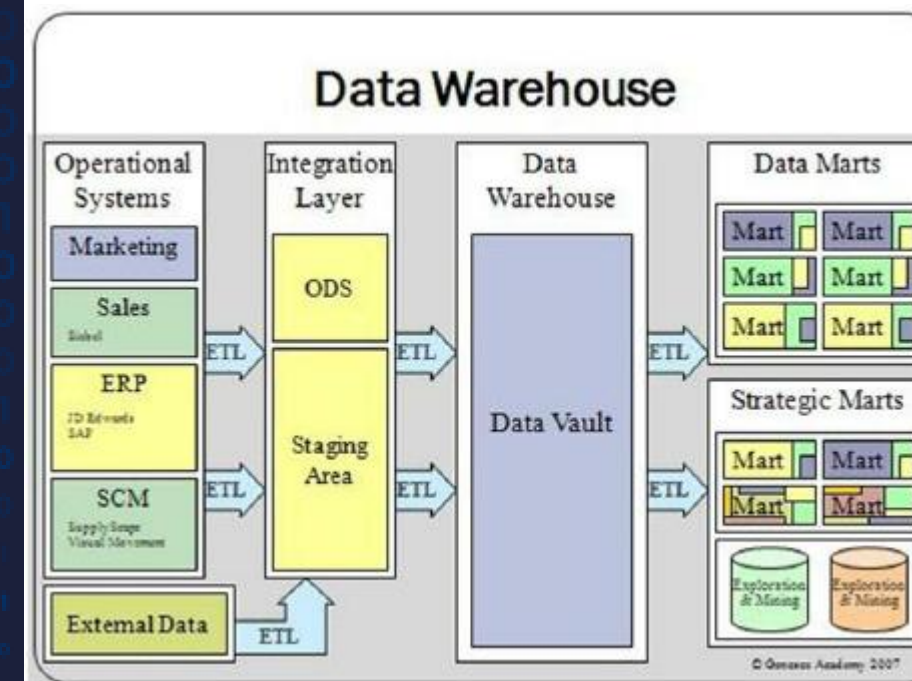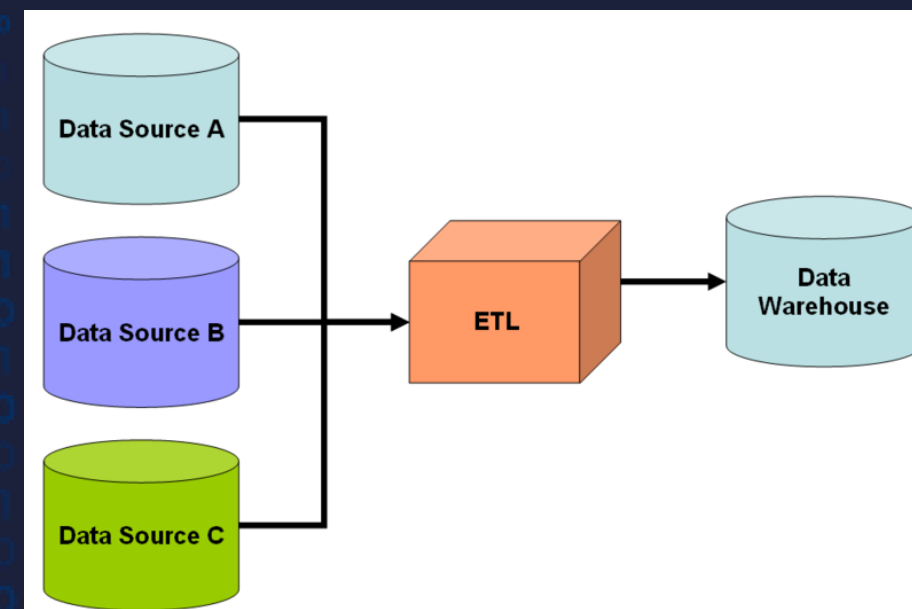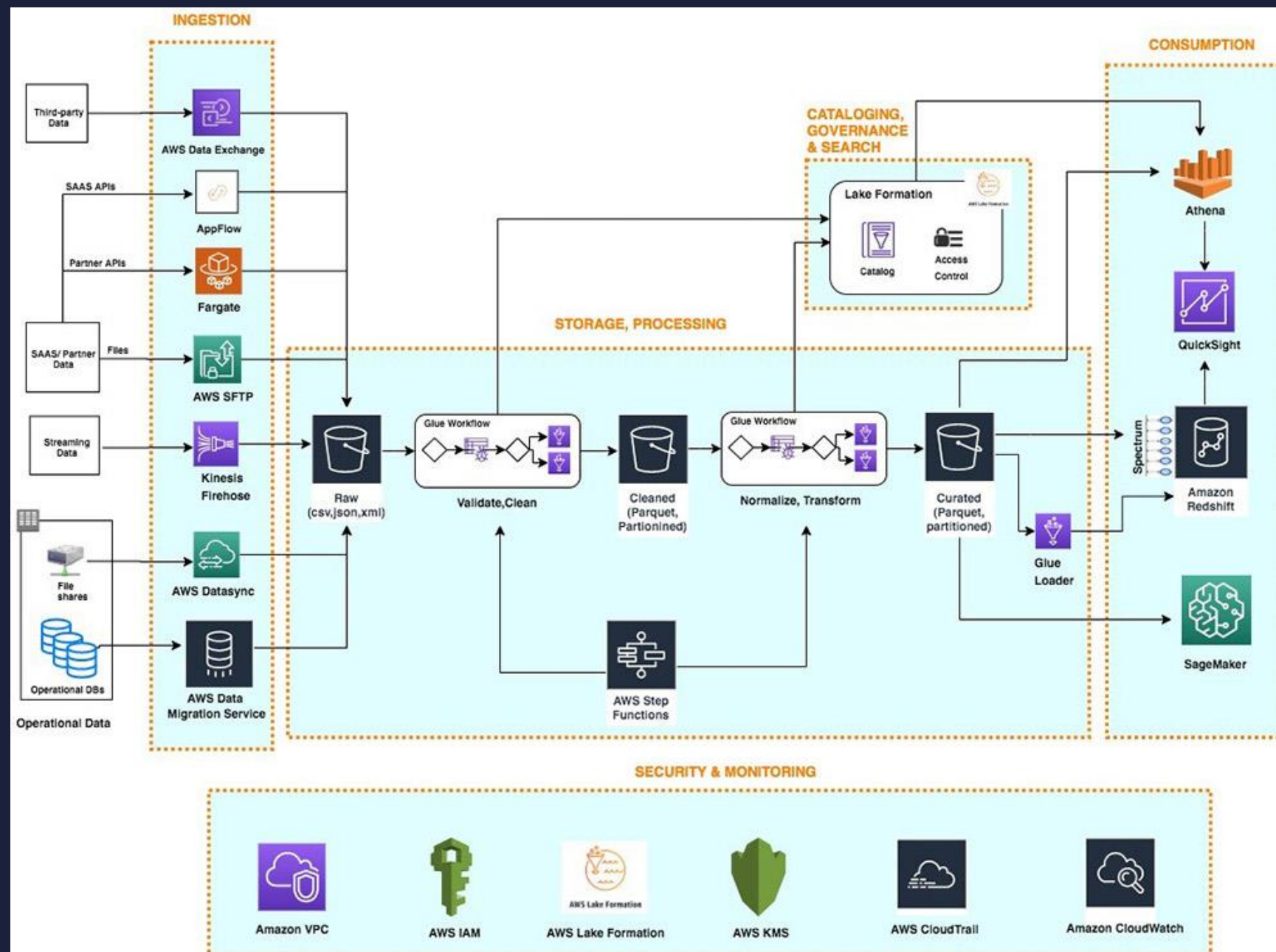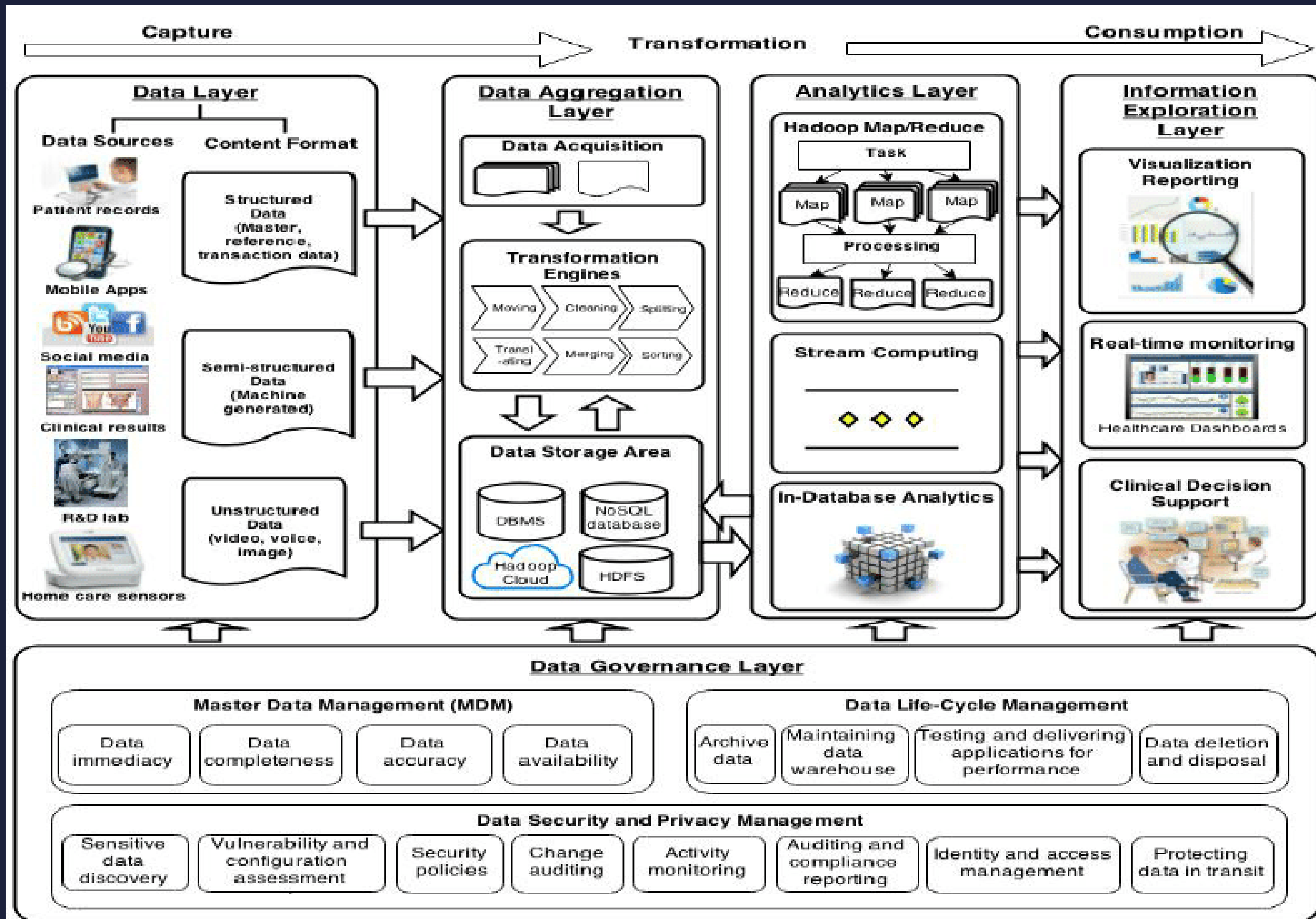
## Supply Chain Optimization:

Big data helps in tracking products throughout the supply chain, optimizing inventory levels, and improving logistics and distribution efficiency.

# Big Data Pipeline

A big data pipeline is a series of processes and tools designed to collect, process, and manage large volumes of data from various sources, transform it into a usable format, and load it into a data storage or analytics system.
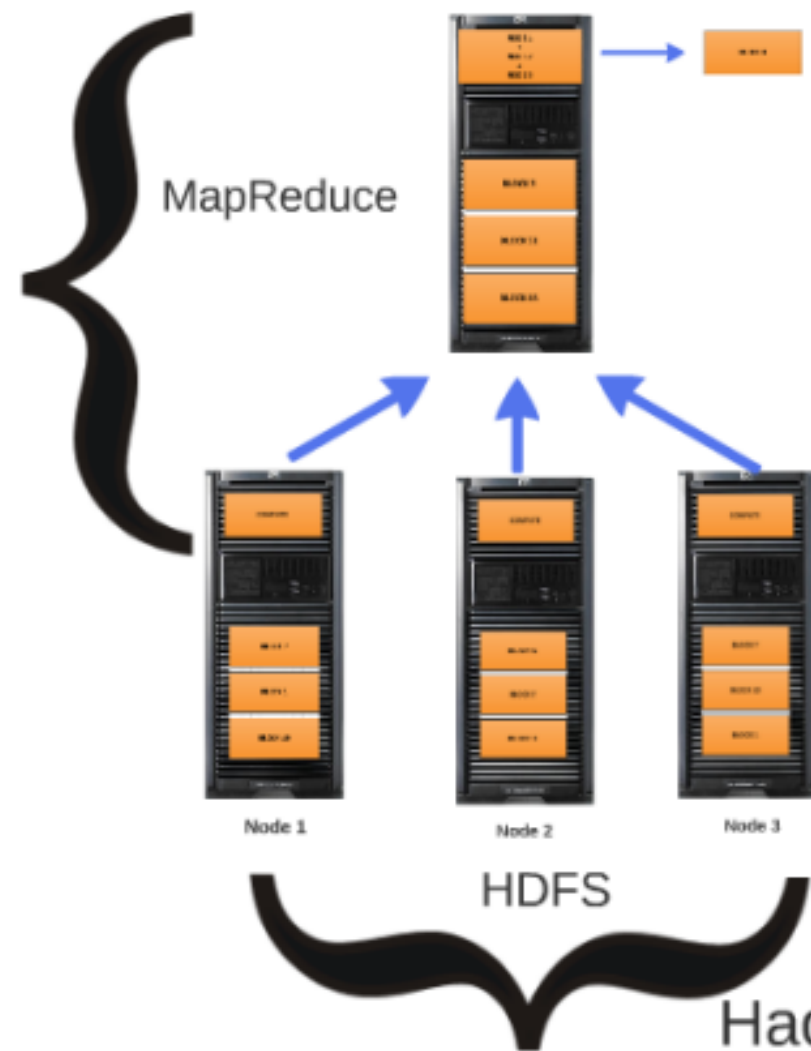
The goal of a big data pipeline is to enable organizations to efficiently and effectively work with massive datasets for analysis, reporting, and decision-making.

# Hadoop Introduction

## PILE OF PAPERS VS. BOOK

VS

Without a file system, information placed in a storage area would be one large body of data with no way to tell where one piece of information stops and the next begins.

# FUNCTIONS OF FILE SYSTEM

- Control how data is stored and retrieved

- Metadata about the files and folders

- Permissions and security

- Manage storage space efficiently

## DIFFERENT FILE SYSTEMS

**Microsoft**

FAT32 - 4 GB File limit 32 GB Volume limit
NTFS - 16 EB File limit 16 EB Volume limit

HFS   - 2 GB File limit 2 TB Volume limit
HFS+ - 8 EB File limit 8 EB Volume limit

**Linux**

ext3 - 2 TB File limit 32 TB Volume limit
ext4 - 16 TB File limit 1 EB Volume limit
XFS - 8 EB File limit 8 EB Volume limit

Why another file system ?

Rack
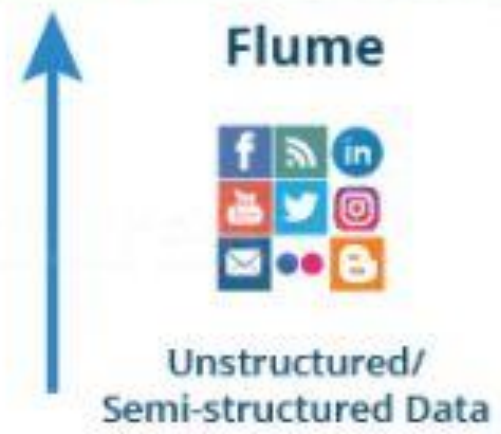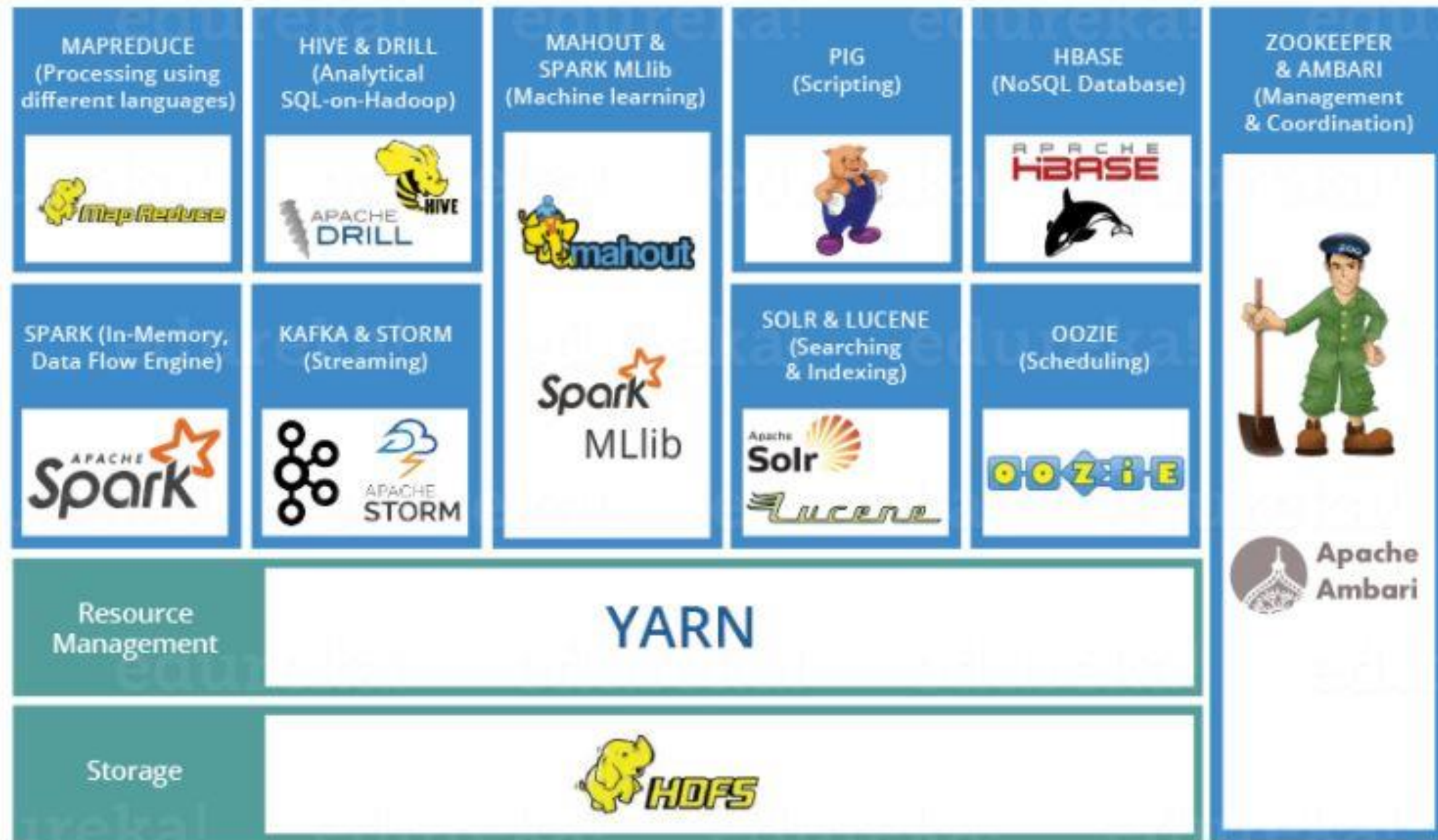
Node

Data Center
(Entire Premises)
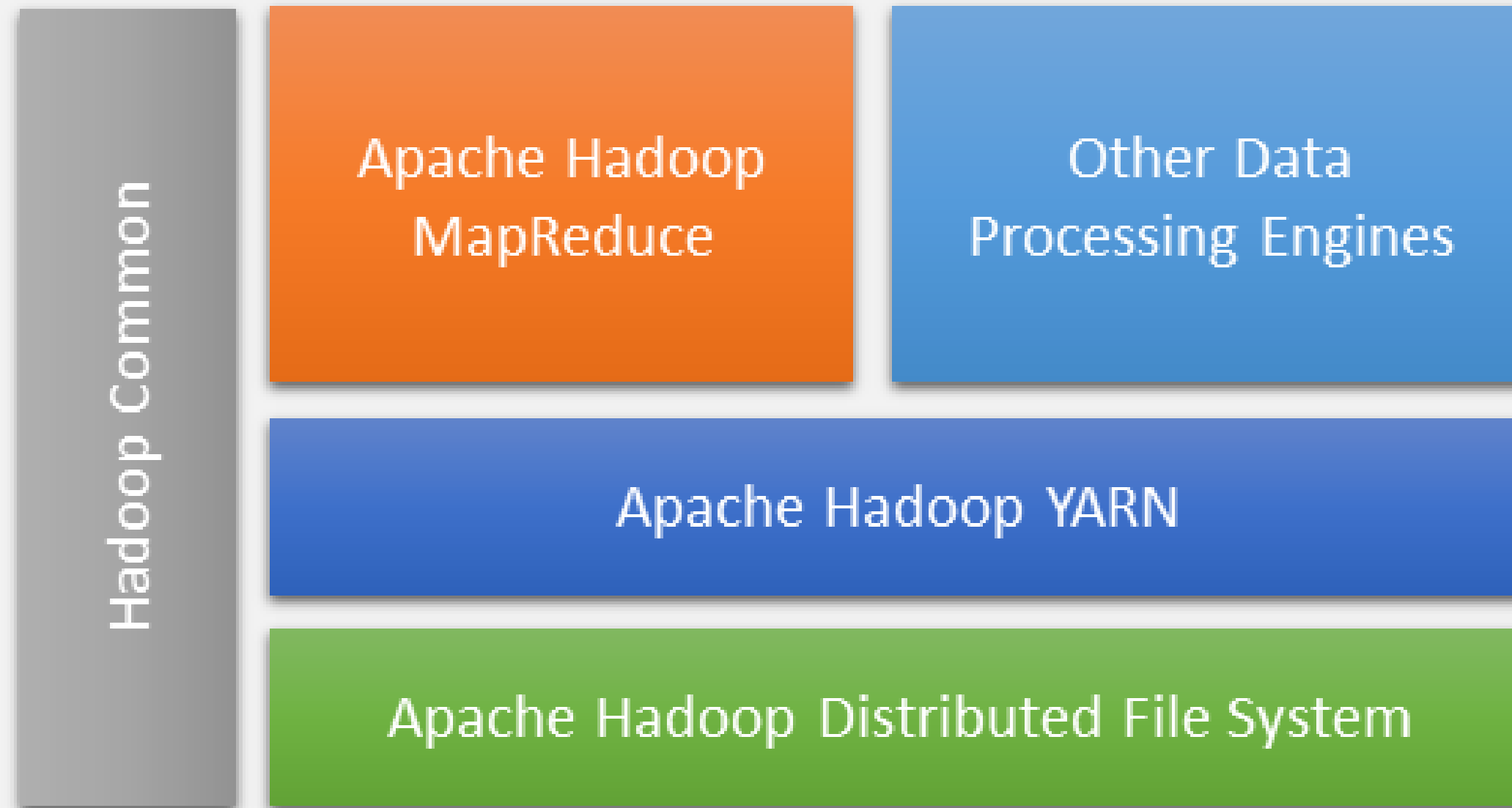
## BENEFITS OF HDFS

- Support distributed processing
  - Blocks (not as whole files)

- Handle failures
  - Replicate blocks

- Scalability
  - Able to support future expansion

- Cost effective
  - Commodity hardware

# Hadoop Architecture

# Write Operation
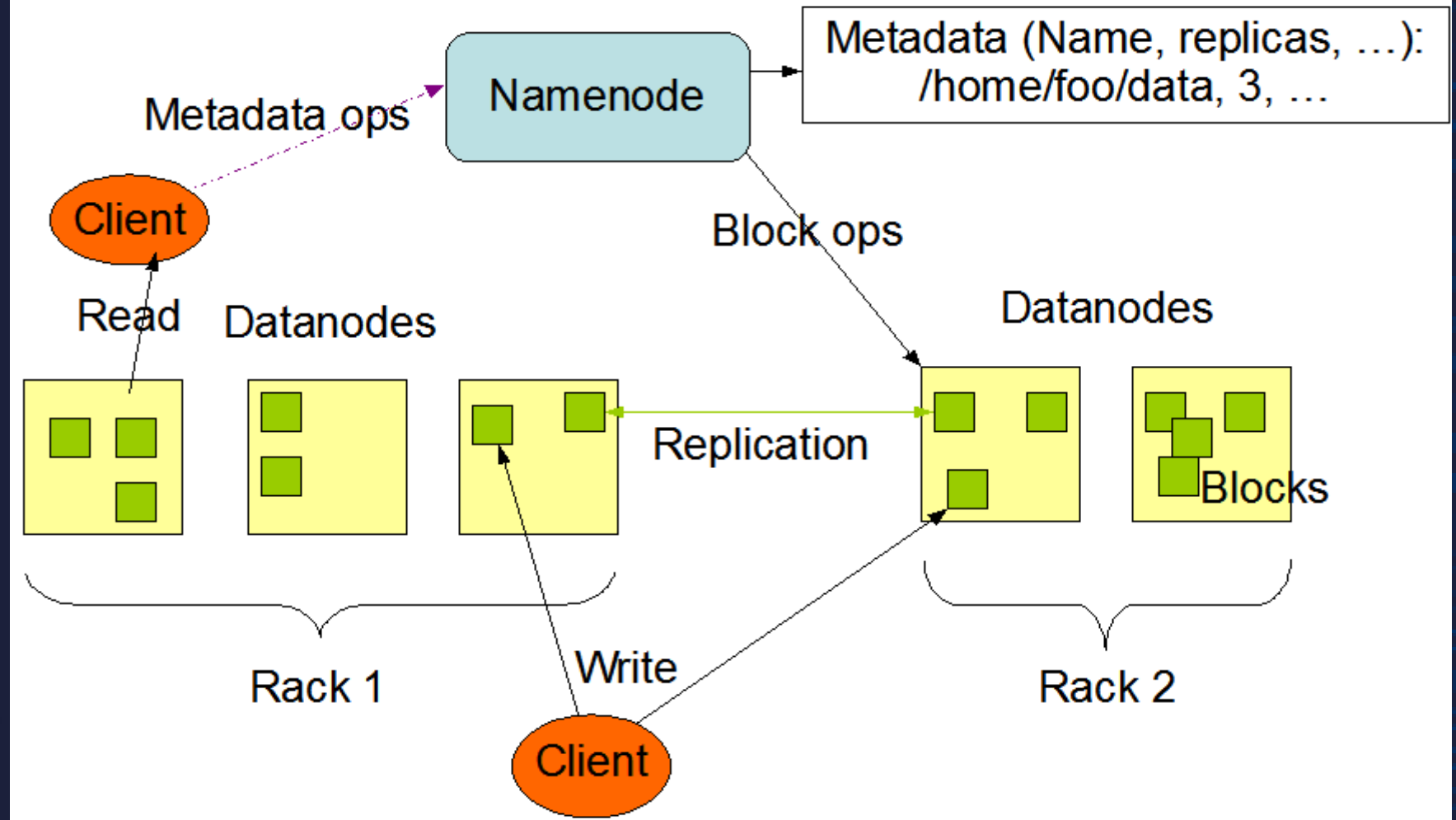
Give me block locations for MySecondFileInHDFS.log

| BLK_0045732 | R8 DN20 | R1 DN2 | R1 DN10 |
| BLK_9610590 | R8 DN20 | R3 DN4 | R3 DN13 |
| BLK_8851205 | R2 DN7 | R1 DN2 | R1 DN10 |

Done

Write BLK_0045732

Done

Write BLK_0045732

Done

Write BLK_0045732

Done

**Name Node**

**Client**

**Data Nodes Pipeline**

**R8 DN20**

**R1 DN2**

**R1 DN10**

# Write Operation - Failure



Change BLK_0045732 to Write BLK_0045732XXX

Give me block locations for MySecondFileInHDFS.log

| BLK_0045732 | R8 DN20 | R1 DN2 | R1 DN10 |
| BLK_9610590 | R8 DN20 | R3 DN4 | R3 DN13 |
| BLK_8851209 | R2 DN7 | R1 DN2 | R1 DN10 |

Write BLK_0045732

(BLK_0045732XXX) Done

Write BLK_0045732

Write BLK_0045732

Done

Write BLK_0045732XXX

(BLK_0045732XXX) Done

**Name Node**

**Client**

**R8 DN20**

**R1 DN2**

**R1 DN10**

**R6 DN12**

**Data Nodes Pipeline**

# Read Operation

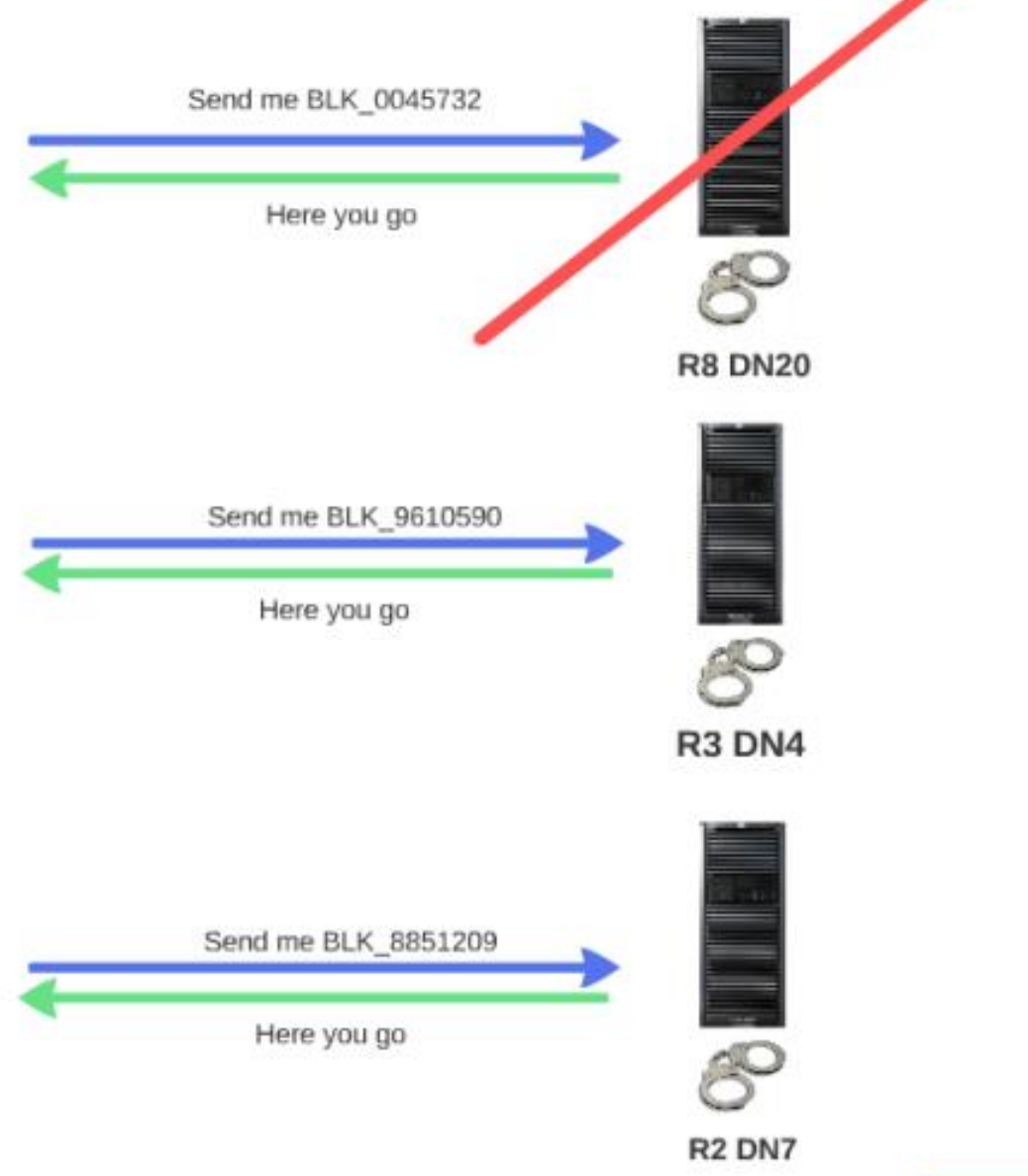**Client**

**Name Node**

**Data Nodes**

Give me block locations for MyFirstFileInHDFS.log

| BLK_0045732 | R8 DN20 | R1 DN2 | R1 DN10 |
| BLK_9610590 | R8 DN20 | R3 DN4 | R3 DN13 |
| BLK_8851209 | R2 DN7 | R1 DN2 | R1 DN10 |

Send me BLK_0045732

Here you go

R8 DN20

Send me BLK_9610590

Here you go

R3 DN4

Send me BLK_8851209

Here you go

R2 DN7