# Big Data and Hadoop

**Venkatesh Vinayakarao**

venkateshv@cmi.ac.in

http://vvtesh.co.in

Chennai Mathematical Institute

Data is the new oil.  - Clive Humby, 2006.

# What Comes Next?

byte

kilobyte

megabyte

gigabyte

??

???

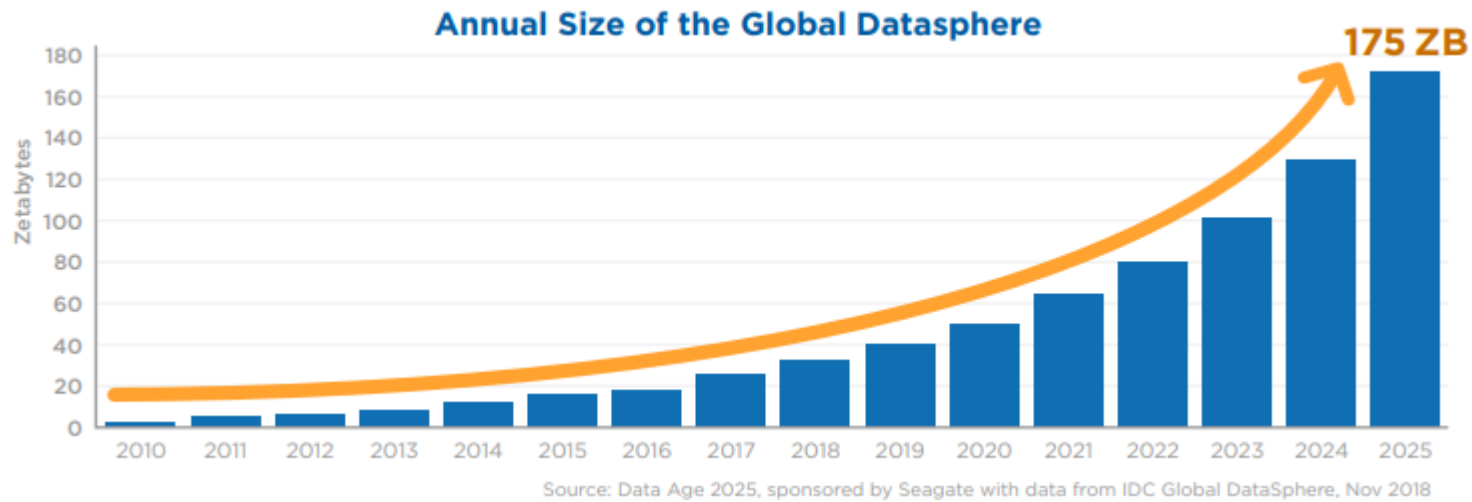????

?????

# Sizes

| Name | Size |
|------|------|
| Byte | 8 bits |
| Kilobyte | 1024 bytes |
| Megabyte | 1024 kilobytes |
| Gigabyte | 1024 megabytes |
| Terabyte | 1024 gigabytes |
| Petabyte | 1024 terabytes |
| Exabyte | 1024 petabytes |
| Zettabyte | 1024 exabytes |
| Yottabyte | 1024 zettabytes |

# Data Growth



**Annual Size of the Global Datasphere**

Source: Data Age 2025, sponsored by Seagate with data from IDC Global DataSphere, Nov 2018
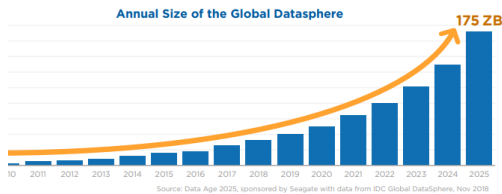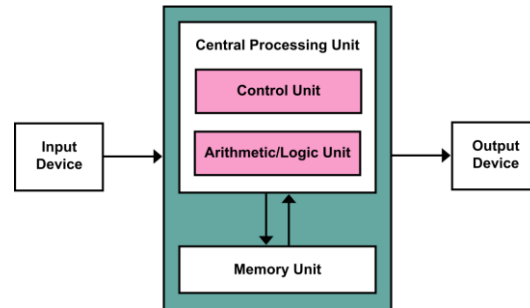
Mankind's quest to digitize the world!
33 ZB (2018) → 175 ZB (2025)
size of global datasphere*

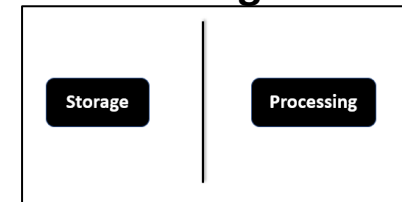*Source: https://www.seagate.com/files/www-content/our-story/trends/files/idc-seagate-dataage-whitepaper.pdf
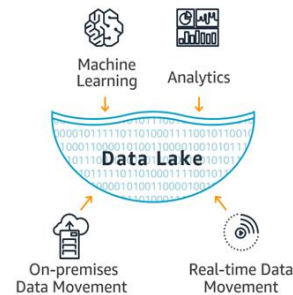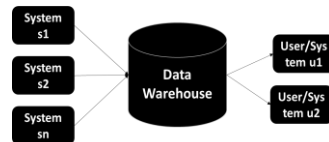
# Evolution of Data and Computers

## Von Neumann Arch

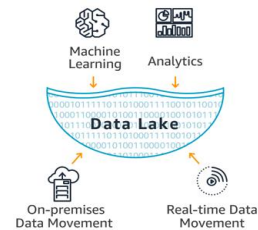**Annual Size of the Global Datasphere**

175 ZB

Source: Data Age 2025, sponsored by Seagate with data from IDC Global DataSphere, Nov 2018

**Central Processing Unit**

Control Unit

Arithmetic/Logic Unit

Input Device

Output Device

Memory Unit

## Challenges

Storage

Processing

## Data Storage

System s1

System s2

System sn

Data Warehouse

User/System u1

User/System u2

Machine Learning

Analytics

Data Lake

On-premises Data Movement

Real-time Data Movement

Amazon S3

**STaaS**

# Recap

## Data Storage



Data Warehouse

Data Lake
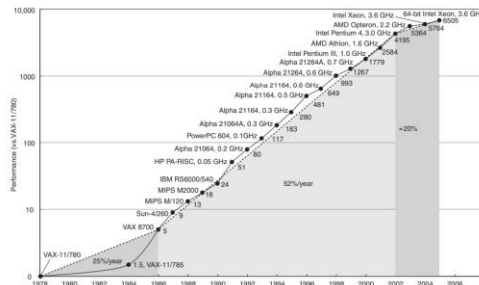
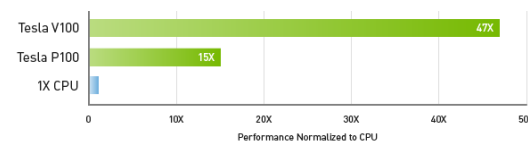Machine Learning

Analytics

On-premises Data Movement

Real-time Data Movement

Amazon S3

STaaS

## Data Processing



CPU Performance

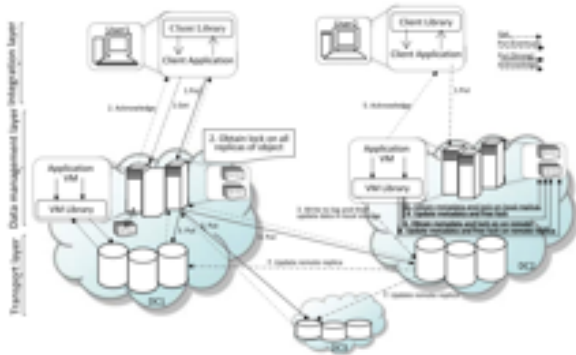47X Higher Throughput Than CPU Server on Deep Learning Inference
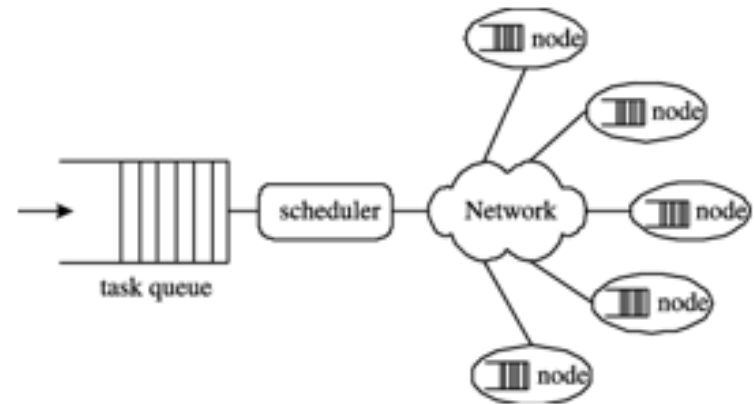
GPU Performance

SuperComputers

# Cloud Computing



Two kinds of Big Data Opportunities

**Storage**

**Processing**

**So, we have the cloud. But, how to store and retrieve data? How to process jobs?**
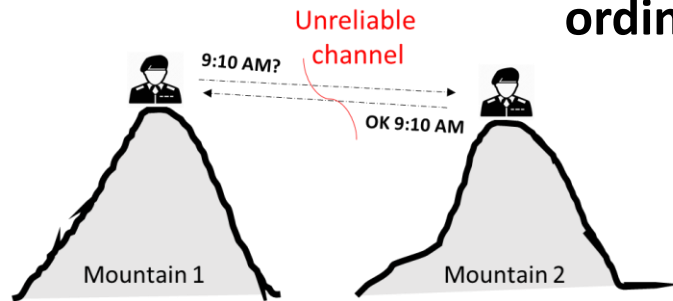
448

# Role of File Systems



**Partitioning**
Multiple OS,
Multiple File
Systems

**Compression**
High Data
Transfer
Time

**Defragmentation**
High Seek
Time

File Allocation,
Free Space
Management
Space
Utilization

Multiple Users

Storage
Media

Multiple Storage Devices

Multi-Tenancy
& data privacy
**Permissions and
Sharing**

Data Variety
**Naming
Convention -
Standards**

Variety of FS
exist
NTFS, FAT, DOS,
CDFS, NFS, …

**File systems are key to handling data.**

449

# Distributed Systems



**Not designed for co-ordination jobs.**

General's Paradox

**WORM Model. Not designed for write-many (interactive) jobs.**

**Not designed for small files.**

# Hadoop and Map Reduce

## When not to use Hadoop?

**No Interactive Jobs**
**No Jobs Requiring Co-ordination**
**No Small Files**

## Hadoop Architecture

Application (map-reduce)
Application (pig)
Application (nosql db)

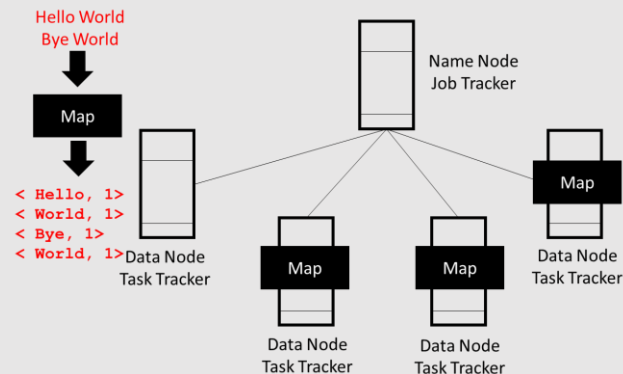**YARN**
(Resource Management – Job Scheduling/Monitoring)

**HDFS**
(Replicated Reliable Storage)

## Map-reduce Model

**Map**

**Shuffle and Sort**

**Reduce**

Hello World
Bye World

Map

< Hello, 1>
< World, 1>
< Bye, 1>
< World, 1>

Data Node
Task Tracker

Name Node
Job Tracker

Map

Map

Map

Data Node
Task Tracker

Data Node
Task Tracker

Data Node
Task Tracker

451

# Map-Reduce Patterns



**Summarization**

**Top 10**

**Filtering**

**Counting**

# NoSQL

Impedance Mismatch
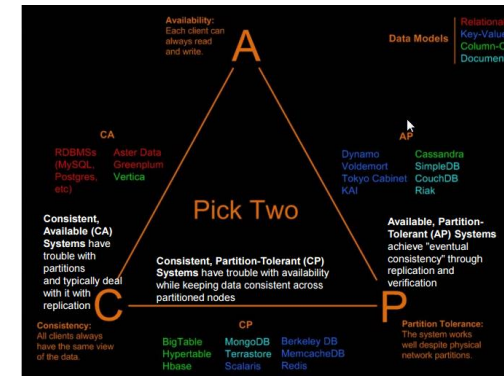


Schema-based Relational Model - maintenance problems
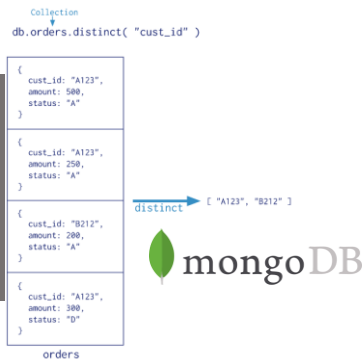
Scale-up Challenges

CAP Theorem

## Types of NoSQL datastores

**Key-Valuecv**
```
redis> GET nonexisting
(nil)
redis> SET mykey "Hello"
"OK"
redis> GET mykey
"Hello"
redis>
```
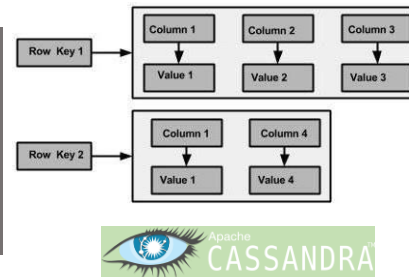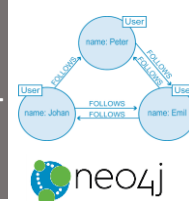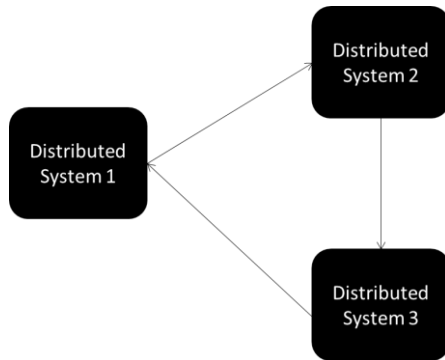redis

**Doc-based**
```
Collection
db.orders.distinct( "cust_id" )

{
  cust_id: "A123",
  amount: 500,
  status: "A"
}
{
  cust_id: "A123",
  amount: 250,
  status: "A"
}
{
  cust_id: "B212",
  amount: 200,
  status: "A"
}
{
  cust_id: "A123",
  amount: 300,
  status: "D"
}

orders
```
distinct [ "A123", "B212" ]
mongoDB

**Columnar DB**
CASSANDRA

**Graph DB**
neo4j

# Web Services

## Interoperability



## CORBA



(client) main()
- Object reference
- Generated stub code
- Object Request Broker

(server) main()
- Object implementation
- Generated skeleton code
- Object Request Broker

network
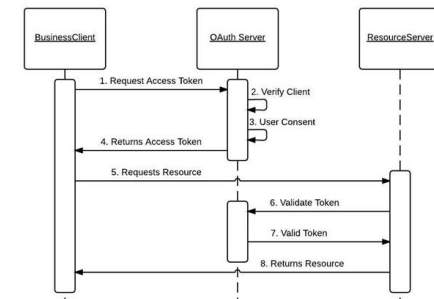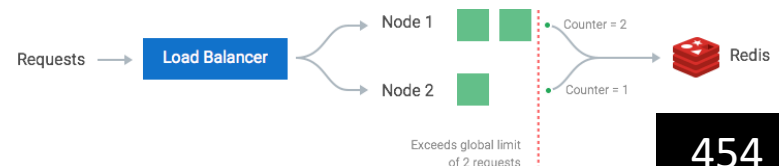
## RMI



## Web Services with REST API



Request a **resource**

Client → Server

Transfer the representation of the state of the resource

## oAuth



## Rate Limiting



## Evolution of Web and App Servers



454

# Building Web Services

# Thank You

Please remember to give elaborate course feedback. I take my course feedback seriously to improve teaching quality including but not limited to the content, presentation materials, and delivery.