| 1980-1990s | 1990s-2000 | 2000s-Till Date |
|---|---|---|
| Mainframes<br>RDBMS<br>MPI<br>PCs | Online applications:<br><br>Email, banking,<br>retail, search | **Big Data and Cloud Computing**<br>Significant explosion of data: |

Distributed scientific computing

Data ware housing:
Online Analytical Processing
(OLAP), BI

Desktop applications moved to web
(Video/photo editing, office tools,
productivity etc.)

SAAS, PAAS, Social and Mobile
became pervasive

Data storage exploded

| No SQL Dbs,<br>MongoDB, Key-value,<br>Column Store | Visualization: D3.js,<br>Kibana |
|---|---|
| Hadoop, MapReduce | SOLR, Lucene |

**Transition from databases to data warehouses to data lakes**
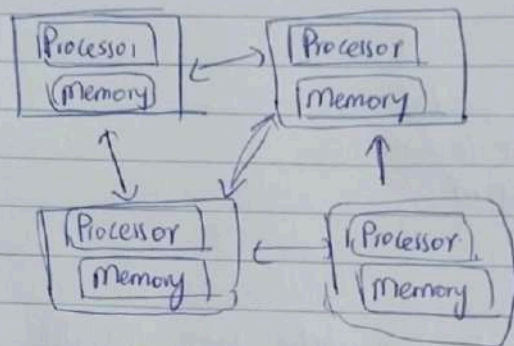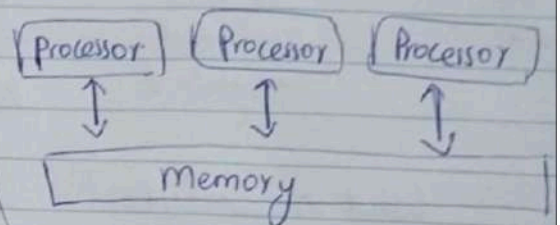
* Curriculum Discussion → * Big Data - Hadoop & spark
                         * Cloud → Azure & GDP
                         → Devop Tools → Git, Github & Doocker.
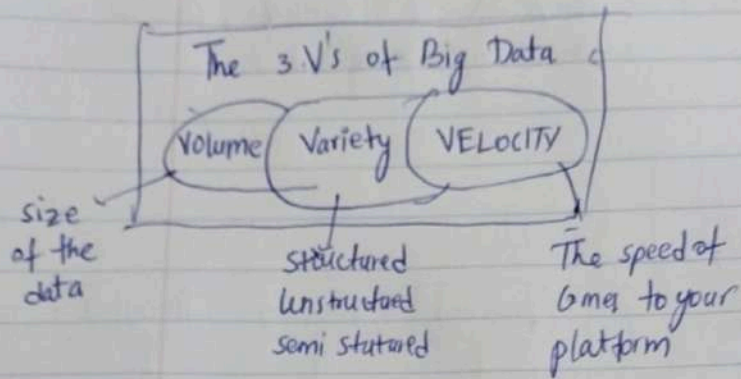                         → No.SQL - mangoDB. Tool.

Processor
Memory
→
Processor
Memory

↑                        ↑
↕                        ↑

Processor
Memory
←
Processor
Memory

Distributed Computing

Processor   Processor   Processor
↕           ↕           ↕
Memory

Parallel-Computing

Cluster Computing
o  A cluster is a type of parallel or distributed proccessing sy
   which consist of a collection of interconnected stand-alone
   coomputers cooperatively working together as a single

**Big data :** It ~~too~~ Data that is Too large & too complex for conventional data tools to capture, store & analyze

The 3 V's of Big Data

( Volume ( Variety ( VELOCITY )

size of the data

Variety:
structured
unstructured
semi stuctured

VELOCITY:
The speed of
time to your
platform

## What is analytics?

The scientific process of transforming data into insight for making better decisions, offering new opportunities for a competative advantage.



① Prescriptive Analytics
enabling smart decisions.
based on data.
what should we do?

Analytics

② Predictive analytics
Predicting the future
based on historical patterns
What could happen?

③ Descriptive aanalytics
Mining data to provide
business insights
What has happanned?

Why do airlane prices change every hour?

Ans ~ Prescriptive analytics
(advice on possible outcomes).

How do grocery cashiers know to hand you coupons you might
actually use?

Ans Predictive Analytics
(understanding the future.).

How does Netflix frequently recommend just the right move?

Descriptive Analytics
(insight into the past).
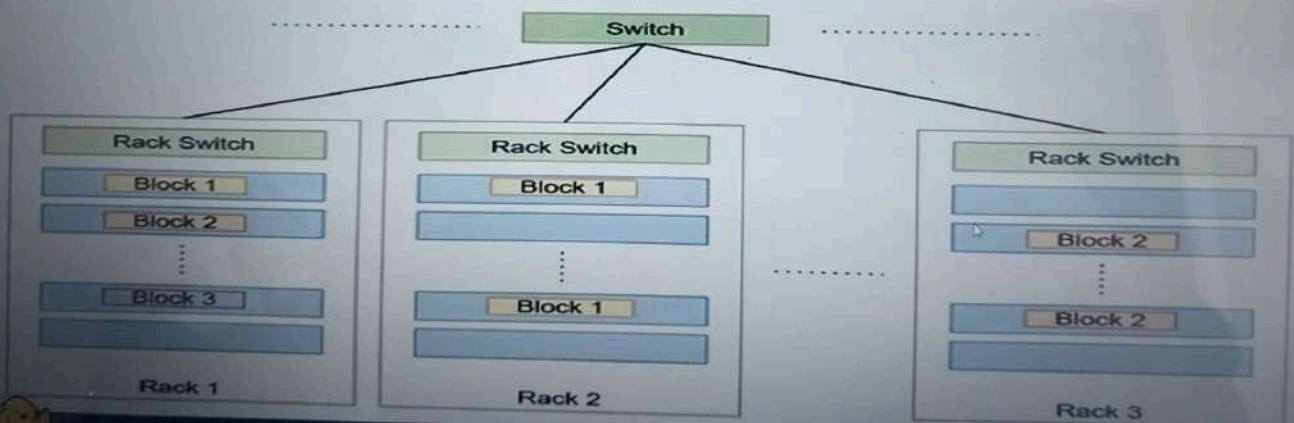
## Problems with Traditional approach

1) Storing huge and exponentially growing datasets.

2). Processing data having complex structure (structured, un-structured,
                                                -semi-structured)

3) Bringing huge amount of data to computation unit becomes a bottleneck

# Rack awareness

Replica storage is a tradeoff between **reliability** and read/write **bandwidth**.

To increase reliability, we need to store block replicas on different racks and Datanodes to increase fault tolerance. While the write bandwidth is lowest when replicas are stored on the same node. Therefore, Hadoop has a default strategy to deal with this **conundrum**, also known as the **Rack Awareness algorithm**.
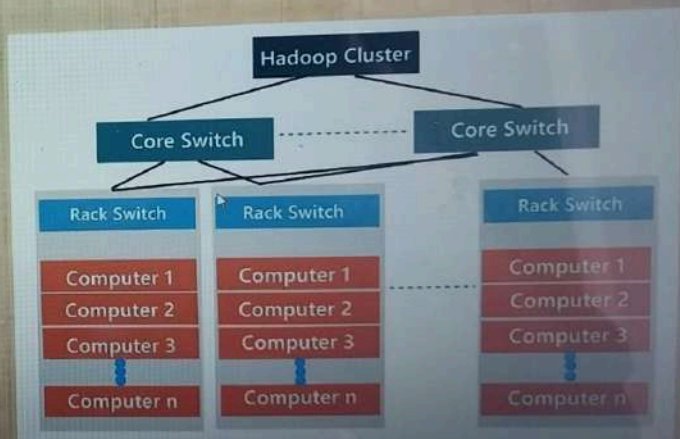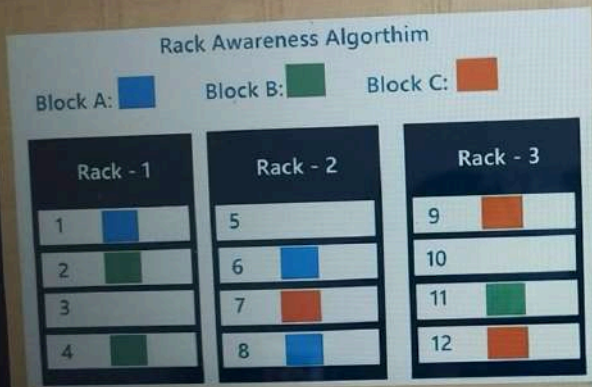
(The data having copies) → by default it has 3 copies

2). Replication →  TA   → 1                           Hadoop → 3.copies

🔲❌ → Data as loost.

A                  A          B              → How replication works)
B                 node1     node2
C        down X    A          A
D                  B          B
                  node 3     node 4           ┌ Data Failure
                                              └→ Network Performance

Rack awarness Algorithm → Replication.

                                        A   B   C   D
                                                        N/w
Node  1 | A   B    10      19                    B 🔲 ↓ 🔲
      2 | C        11      20                      performance
      3 | A   D    12      21                      go-down
      4 |     A    13      22
      5 | B        14      23
      6 |          15/17   24
      7 |     B    16      25
      8 | C   D    17      26
      9 |     D    18      27

      Rack1        Rack2      Rack3

         ↳    The performance is high ✓ but it we loit Rack1 we can
                                          loose total data so.
                                            it will store.
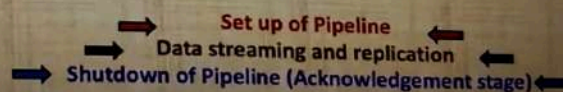                                            atmoit 2 Rack

      Rock Awarness algo → at atmoit -2 Rocks.
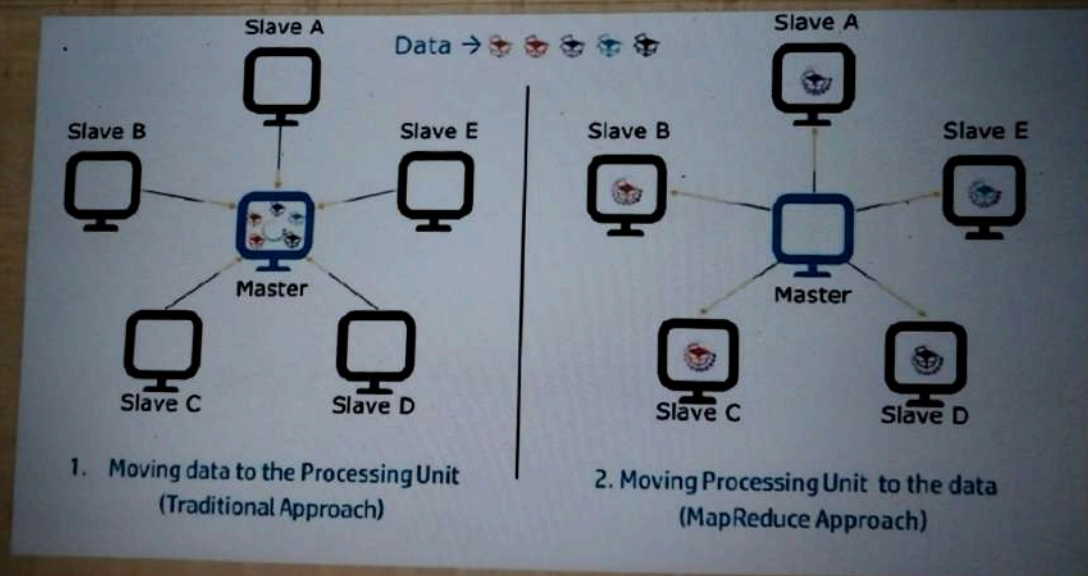
# HDFS Read Architecture:



HDFS - Read Architecture

# Advantages of MapReduce

## 1. Parallel Processing:



1. Moving data to the Processing Unit (Traditional Approach)

2. Moving Processing Unit to the data (MapReduce Approach)

# Map Reduce:

**T.A** ~~Total~~ (Tradilational approch)  |  **B.D.A** (Big Data approch)

(erver)

Accessing

① In the processing of trafer the data one machine to another machine. anthing can happenned.

```
1 A
123
4 5 6
```
h   5

3

elient set data

② data
   1 server

one single server process
the entire data.

---

SlaveC
```
1
```

SlaveA
```
2
```

Slave D
```
5
```
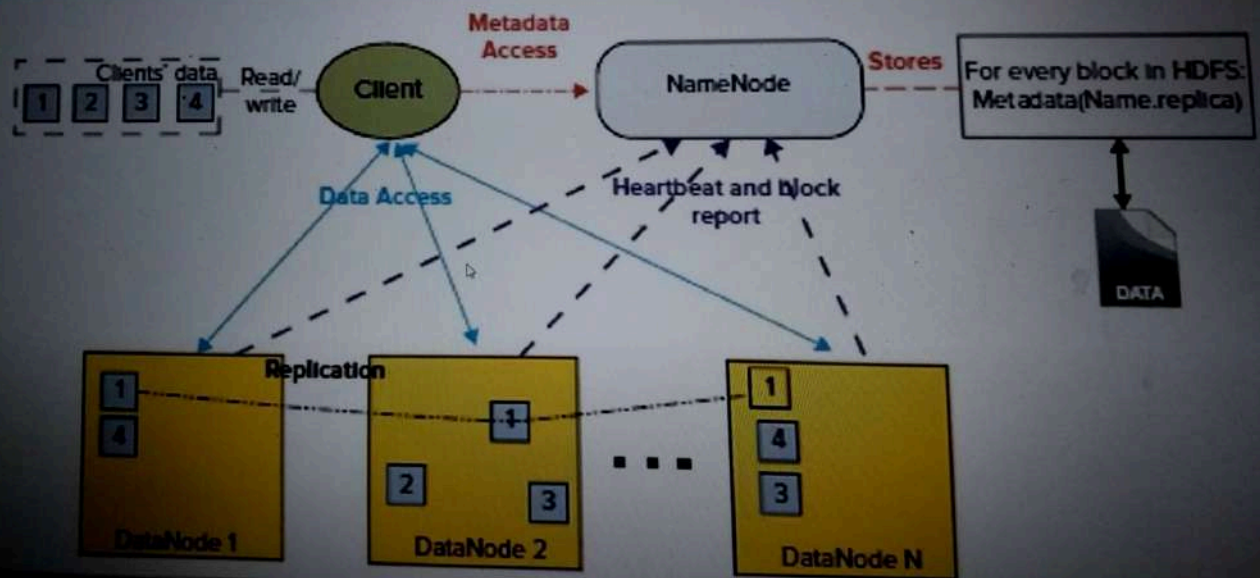
master

SlaveB
```
3
```

Slave E
```
4
```

① here the data
Sloved) won't tranter to one machine to another.

② The processing will take
sloved) care of each slave.

③ It will only share the results of the slave's

## core components or daemons:-

## Hadoop Architecture:-

# Hadoop 1.X Limitation

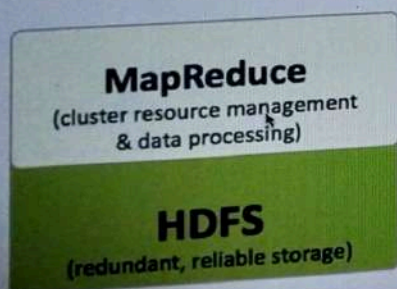Main drawback of Hadoop 1.X is that MapReduce Component in it's Architecture.
- → If the Name node server is down than total will down
  - → It not suitable for Real-time Data Processing / Data Streaming
  - → Job Tracker is the single point of failure
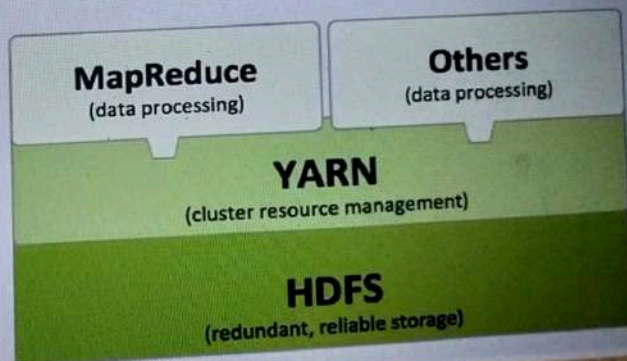  - → it runs only Map / Reduce jobs

# Hadoop 2.X
### YARN (Yet Another Resource Negotiator)
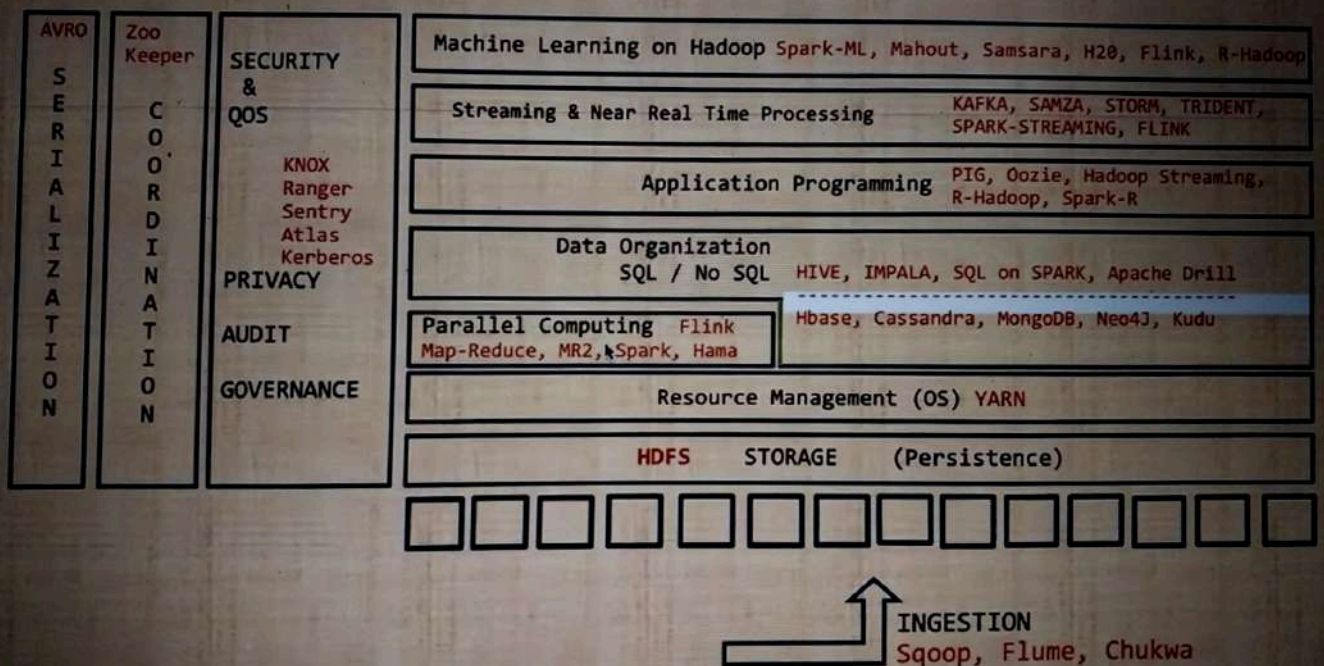
# Hadoop 2.X Allows to work in MR as well as other distributed Computing models like spark, Mama, Giraph, Message Parsi

# 2.X Has better scalability, Scalable up to 10000 nodes per clu

# Multiple Namenode servers manage multiple namespace

# Can serve as a platform for a wide variety of data analytics - possible to run event processing, streaming and real-time operations.

| AVRO | Zoo Keeper | SECURITY & QOS | Machine Learning on Hadoop Spark-ML, Mahout, Samsara, H20, Flink, R-Hadoop |
|---|---|---|---|

**AVRO**

**S E R I A L I Z A T I O N**

**Zoo Keeper**

**C O O R D I N A T I O N**

**SECURITY & QOS**

KNOX
Ranger
Sentry
Atlas
Kerberos

**PRIVACY**

**AUDIT**

**GOVERNANCE**

Machine Learning on Hadoop Spark-ML, Mahout, Samsara, H20, Flink, R-Hadoop

Streaming & Near Real Time Processing    KAFKA, SAMZA, STORM, TRIDENT, SPARK-STREAMING, FLINK

Application Programming    PIG, Oozie, Hadoop Streaming, R-Hadoop, Spark-R

Data Organization
SQL / No SQL    HIVE, IMPALA, SQL on SPARK, Apache Drill
--------------------------------------------------
Hbase, Cassandra, MongoDB, Neo4J, Kudu

Parallel Computing    Flink
Map-Reduce, MR2, Spark, Hama

Resource Management (OS) YARN

HDFS    STORAGE    (Persistence)

INGESTION
Sqoop, Flume, Chukwa

— 1.X — Job tracker
  └ Resource managment & processing
2.X.V          YARN                        ⊃SPARK

Hadoop 2.X Version
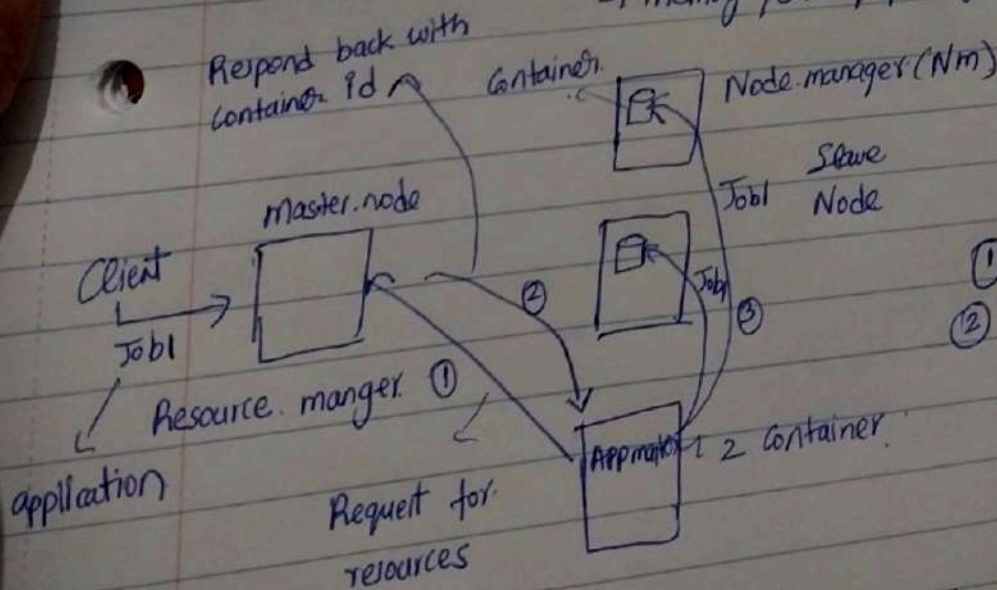    ├→ YARN  → Resource management
    └ SPARK  → Processing

YARN → Yet Another Resource Negotion

→ Container Concepts
    └ memory /duck /cpu. cycles         Job:app.marter
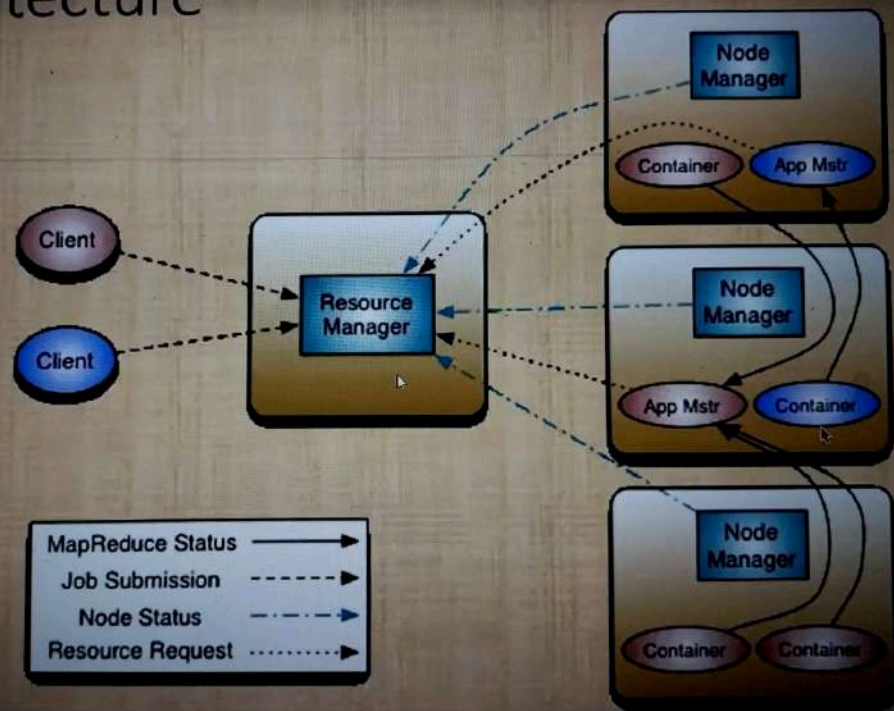                                        1:1 sure.

Respond back with            Container.      Node.manager (Nm)
container id ^                                              Slave
                                              Jobl    Node          Resoure manager.
         master.node                                                ①.It creater a container
Client →  [      ]                          [    ]                  ② It create app marter
   └                                          Job ③
   Jobl
   / Resource. manger ①                    [APP marter] 1 2 Container.
application      Request for.
                 resources

The node manager will give heart beat to the Resource manger.
health of each Node Is updated to node manger to resoure ma

# Running a YARN application
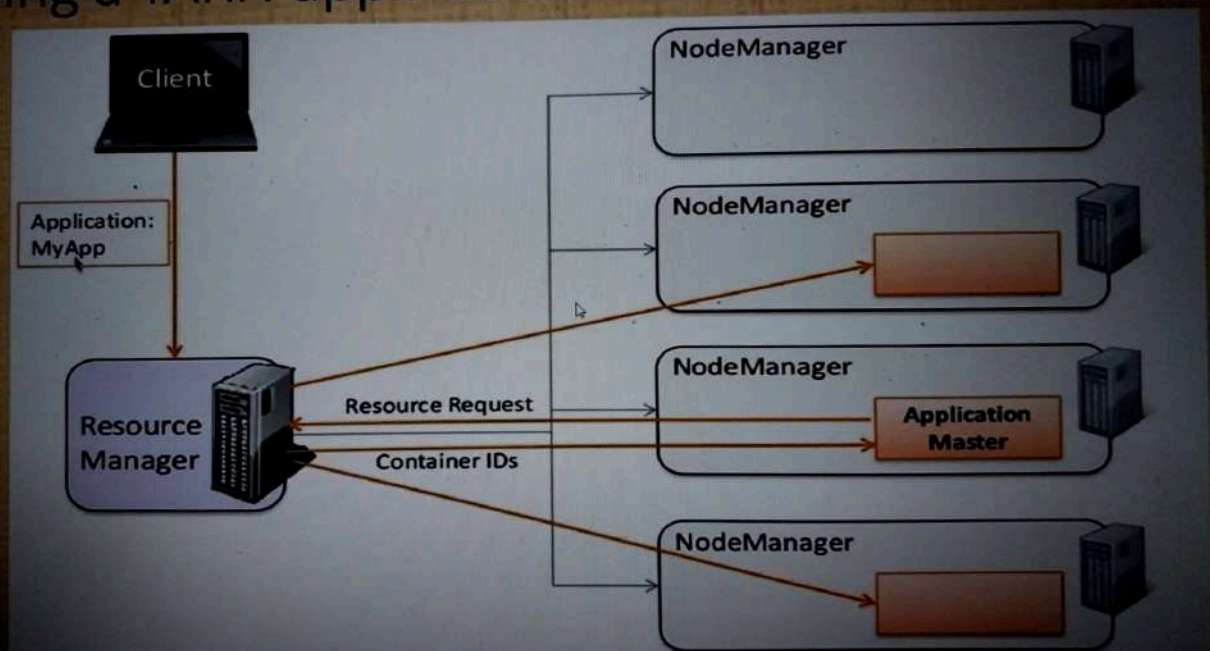
Running Apps on YARN (2)

# CapEx vs OpEx

| Capital Expenditure | Operational Expenditure |
|---|---|
| Spending on infrastructure is completed upfront | No upfront cost |
| Cost written off over a period of time | Pay for service as you consume it |
| | Deduct from tax bill in same year as expense occurs |

# Typical On-premises CapEx Costs

Server costs

Storage costs

Network costs

Backup and Archive Costs

Datacenter Costs including DR

# Why move to the cloud?

## Current

- Focus on building and deploying applications
- Maintenance is done for you
  - No more software patching, hardware setup, upgrades and IT management

## Reliable

- Your data is safe
- Cloud vendor provides:
  - Data backups
  - Disaster recovery
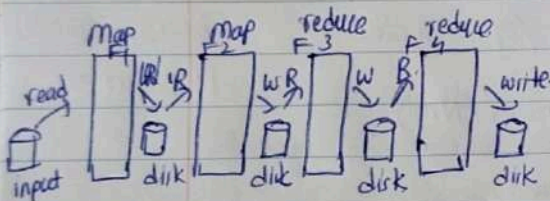  - Data replication

Spark → Parallel processing platforms

(In-memory computation)
↑
Spark

M.R (map Reduce)

→ 2 Functions map & Reduce
→ batch processing
→ Disk to process the data

→ multiple Fuctions
→ both batch & real-time data
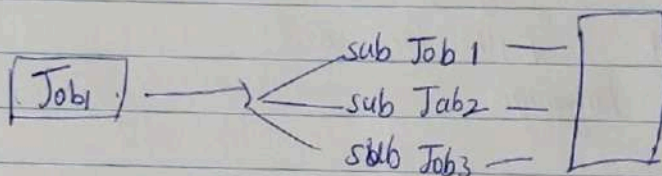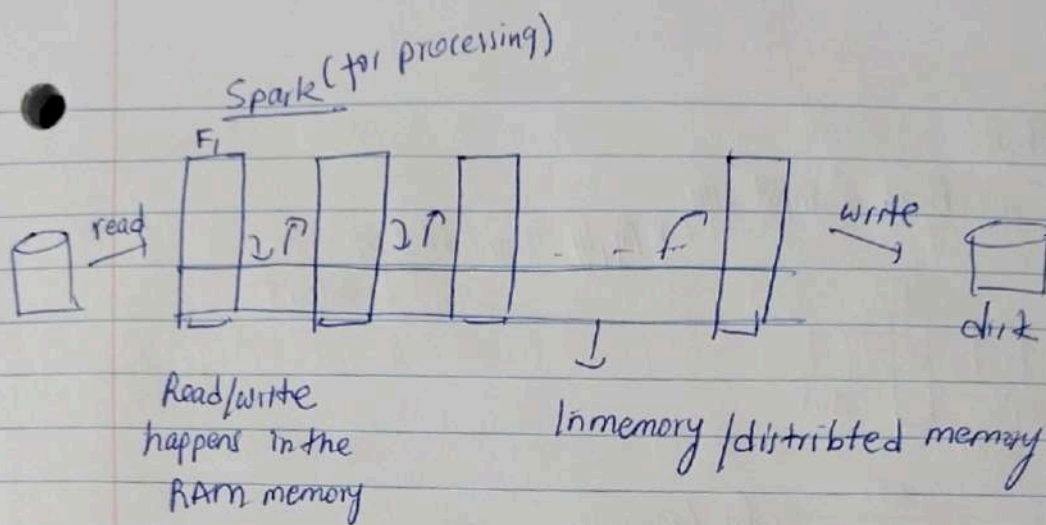→ Memory → process the data



→ Multiple Read/write operation
→ Time - processing is more/high

on Disk
takes more
time

RAM → Memory → we can reduce time

ROM → Disk → Time
                is huge

Spark (for processing)

F₁

read → [2 ↑] [2 ↑] [ - - f ] write → disk

Read/write
happens in the
RAM memory

Inmemory/distribted memory

[Job₁] → sub Job 1 ——
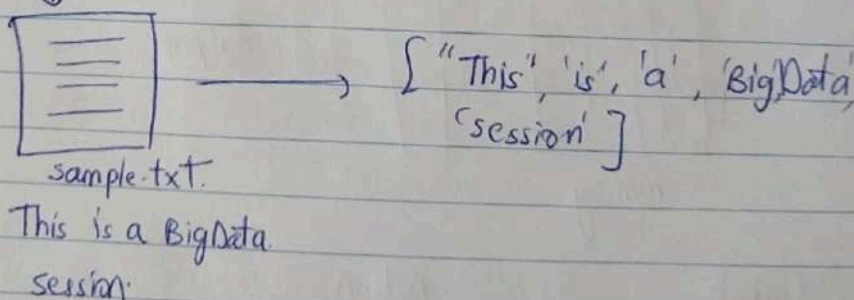          sub Jab₂ ——
          sblb Job₃ ——

Spark

2 operation        Converting one form of data into another form
① Trasformation —— Multiple Fuction
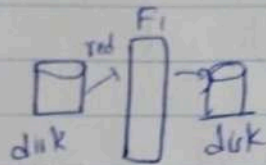② Action → Multiple Fuction
        ↳show the result

[≡≡≡] ————→ [ "This", 'is', 'a', 'BigData',
sample.txt.              'session' ]

This is a BigData
session.

## Python

Read a input file
&#8627; Apply Transformation
&#8627; output.



disk — Trsf → F₁ — disk

Spark ——writen—— Scala Language
+ &#76; pyspark = python language
Python &#8627; language

## Py-Spark

lazy Evaluation { Transformation - unique feature
&#8759; will not consume any
cpu cycle / ~~execute~~ not execute until.
we call action



F₁    F₂    Action

→ output.

memory

| Transf | Action |
|---|---|
| • map() | • collect() |
| • filter() | • show() } &#9758;'shows output |
| • Reduce By Key() | • take() |

spark —: API's

    spark → 3.X.V

    1.4 ≤ 1.4 >
     ↓      ↓
    RDD API  DF ← RDD
                  & 
                SQL

                                 output
RDD —: unstructured Data  —  [ ....... ]
 ↓
DF —: structured way of processing Data. —
&SQL
      L Hand-on — cloud

             master.    Slave.
HDFS  —  Name Node  Data Node.
MR  —  Job Traker  Task Traker
YARN  —  Resource manager  Node manager.
SPARK —  Drivers       Executors.

use case - clould? why I should move to cloud

start up - in developing a mobile App - Food ordering . App

End - End product Development

client → Development → Deliver (customer)
                                    launch it
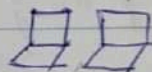
Infrastructure:     →  server            Benglure - 10,000
(server, computers)     OS                            customer
                        storage
                        software stack < tool
                                        funcality
                                ↓              → Itrastructure for.
                        Avg computing            production (launch it)
                            power.               is different compare to
                        🗇🗇                      development.
                        p Laptops.
                        10-developers

Day 1        Target.            Intra required
             customer           r,000    s,000.
             ―――――              ┌──┐    ┌──┐        IT. Team
             10,000             │  │    │  │
                                │  │    │  │        2. members
                                └──┘    └──┘
                                server1  server2

             No. of customer - 2k - 5k

Day 5    -     9k customer → works fine

Day 6  -       15k Customer → Accell the App
                   ? - server down
                     - customer - bore & bad Feedback

Add New servers → 2 server – Budget
                                    – configure.

Day 10       20,000                        4 server's    – 3 IT
                No issuers. ✓                                     Team.

Day 20. → 30,000                           6. server  – 4 IT
      |         ✓ 10 day – No. issues ✓                                Team

Day 30.
   Day 32 –        No. of user's – 5K           Cost is
    33                        "       – 6K    ?    fixed
    "                        "      –10K            ?

observation

Focus            ⎡ → maintance isues – costly , Time
   ↑                 | → server cost
                | → IT. Team.
Developing      ⎣ → Cost is fixed
New product
with
less budget.        3rd party — comes to you
               (They take of' all
               Infastructure (serve's)   after — cloud.
                Accepted                 ↑
                                       pay for what
                                         you use.
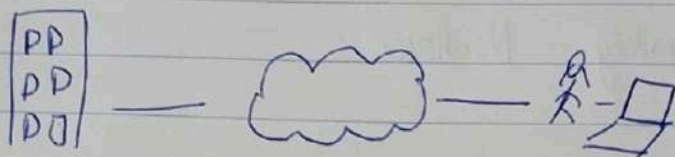                        Day – 4 hrs – pay for 4 hrs
                          1hr → " "

Where can infrastructure be hosted?
• For an enterprise, it is on a Data Centers

Whath is cloud Computing
• Delivery of computing services over internet



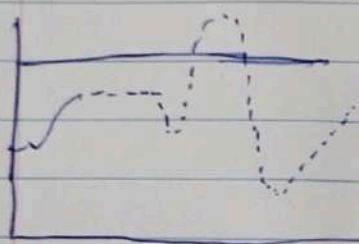—) I aaccessing someone else over computer over the internet for my computer and pay only for cloud services you use.

What is Cloud Computing?
* Renting resources vs Purchasing the hardware
# Pay for what you use
* Run your applications in someone else's datacenter
* Cloud provider is responsible for the physical hardware and facilities necessary to execute your work
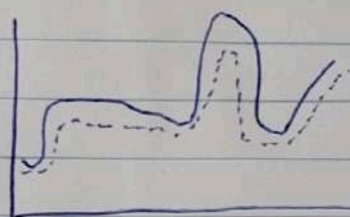* Cloud provider is responsible for keeping the server to up todc

Elasticity)

→ As your workload changes, resources can be changed to compensate
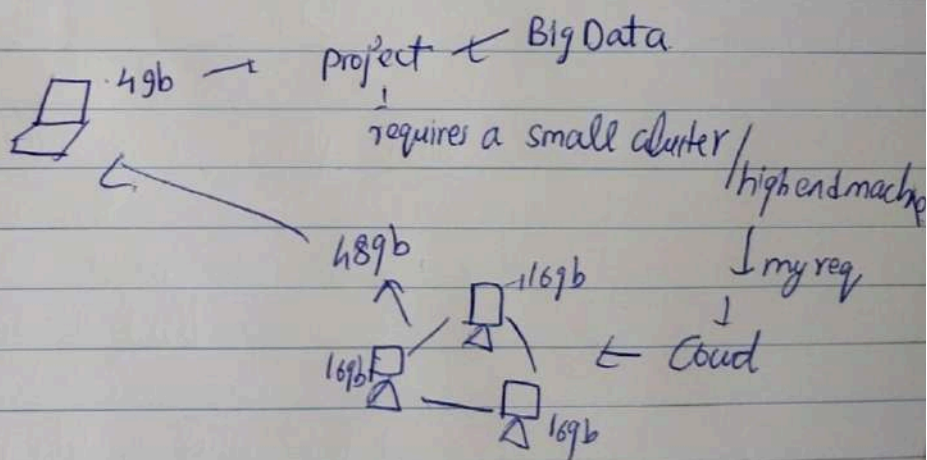example : Seasonal Demand for retail website Black Friday (up or down).

resources (server)
——→)

←load)



static



Elastic Scaling



.4gb → project ← Big Data
↓
requires a small cluster /high end machine

4.8gb
↑
16gb
16gb          ↓ my req
↓
← cloud
16gb

16gb

# Pizza as a Service

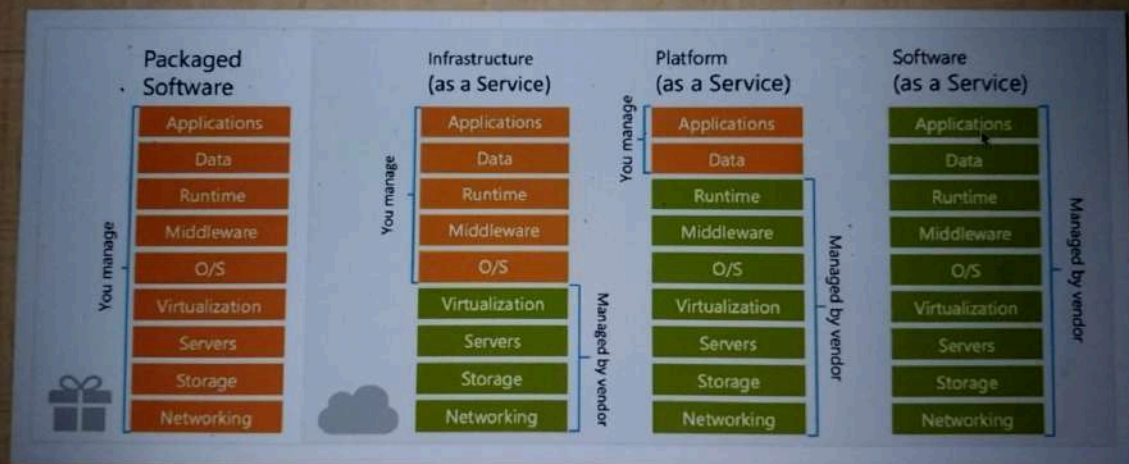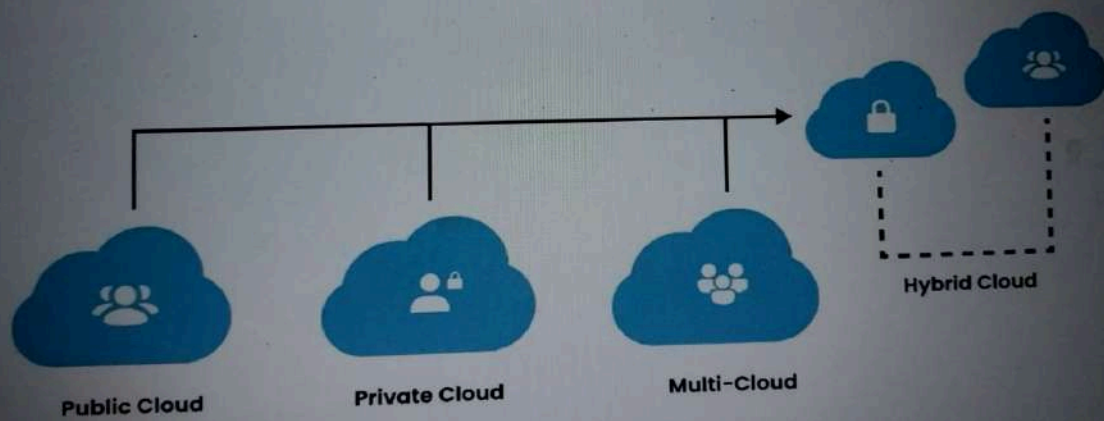| Traditional On-Premises (Legacy) | Infrastructure as a service (IaaS) | Platform as a service (Paas) | Software as a service (Saas) |
|---|---|---|---|
| Dining Table | Dining Table | Dining Table | Dining Table |
| Drinks | Drinks | Drinks | Drinks |
| Electric / Gas | Electric / Gas | Electric / Gas | Electric / Gas |
| Oven | Oven | Oven | Oven |
| Fire | Fire | Fire | Fire |
| Pizza Dough | Pizza Dough | Pizza Dough | Pizza Dough |
| Tomato Sauce | Tomato Sauce | Tomato Sauce | Tomato Sauce |
| Toppings | Toppings | Toppings | Toppings |
| Cheese | Cheese | Cheese | Cheese |
| **Made at Home** | **Take and Bake** | **Pizza Delivery** | **Dining Out** |

■ You Manage   ■ Vendor Manages
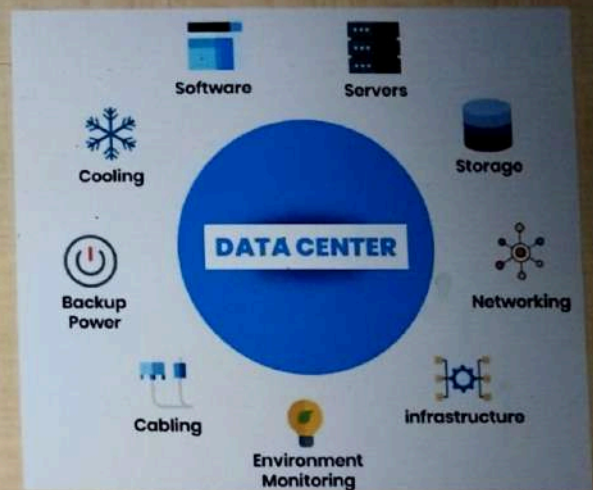
# XaaS (What can be rented?)

Types Of Cloud Deployment Models

# What is Azure?

- A Cloud Computing Platform from Microsoft

- Released as **Windows Azure** in February 2010
  Renamed to **Microsoft Azure** on March 25, 2014

- Provides a web portal to access and manage cloud services & resources.

- Free to start, pay-per-use

# Azure Global Infrastructure

1. Data centers
2. Regions
3. Geographies
4. Availability Zones
5. Region Pairs

# Azure Regions

Location for your resources

Area containing at least one datacenter

Select a region when deploying a resource

# Azure Geographies

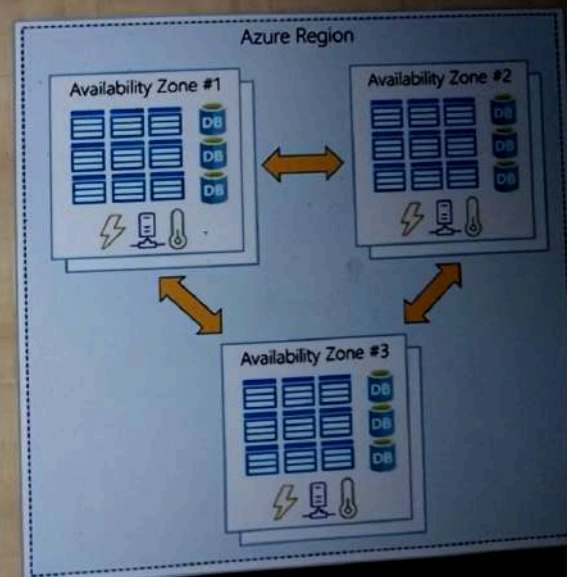An Azure geography is an area of the world that contains at least one Azure region.

Ex: United States, United Kingdom, India, Asia Pacific etc
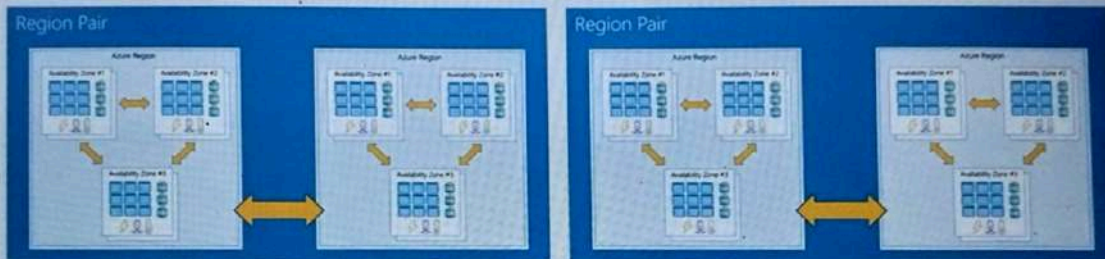
## Azure Government

- This geography is only available to the United States federal, state, local, and tribal governments and their partners.

# Azure availability zones

Availability zones are physically separate datacenters within an Azure region. Each availability zone is made up of one or more datacenters equipped with independent power, cooling, and networking. An availability zone is set up to be an *isolation boundary*. If one zone goes down, the other continues working. Availability zones are connected through high-speed, private fiber-optic networks.

Each Azure region is always paired with another region within the same geography (such as US, Europe, or Asia) at least 300 miles away. This approach allows for the replication of resources (such as VM storage) across a geography that helps reduce the likelihood of interruptions because of events such as natural disasters, civil unrest, power outages, or physical network outages that affect both regions at once. If a region in a pair was affected by a natural disaster, for instance, services would automatically failover to the other region in its region pair.
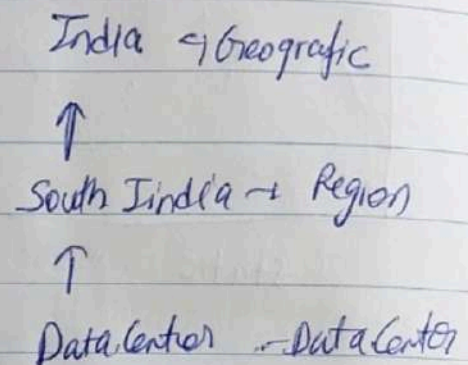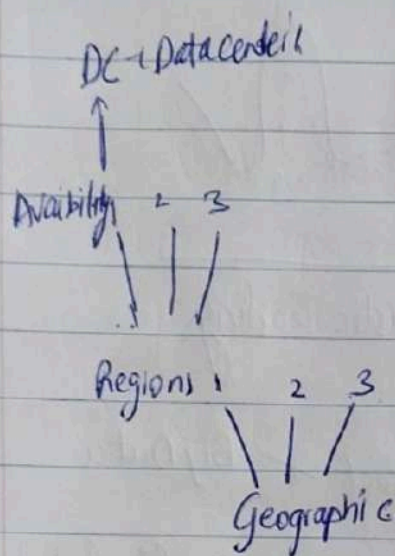
Examples of region pairs in Azure are West US paired with East US and SouthEast Asia paired with East Asia.

# Azure region pairs



| Geography | Regional Pair A | Regional Pair B |
|---|---|---|
| Canada | Canada Central | Canada East |
| China | China North | China East |
| India | Central India | South India |
| Japan | Japan East | Japan West |
| North America | East US | West US |

# Core cloud services

Compute      Storage      Networking
       App Services      Analytics

DC → Data centers
↑
Availability   2   3
    \ | /

Regions 1    2   3
       \ | /
       Geographic

India → Geografic
↑
South Iindia → Region
↑
Data Center → Data Center

## IAAI

- → Compute Engine
- → Storage servers
- → Network services

Compute engin       Storage servic       network server
     ↓            ↓

→ VM        Structure data. → Azure SQL

→ Containers     Semi-Structure data → Cosmas DB    → NIC

→ Kubernetes    un-structured data → ADL &     → IP. address
                         Blob      → subnets

— Resource Group

R.G ......                    Service — Azure
 ├ Vm
 └ storage

# Storage Services.                    ⌒ col / attributes
(store any data)                 rows⌐ ▦
 ├ Structured Data — Azure  tuple⌐      ← schema.
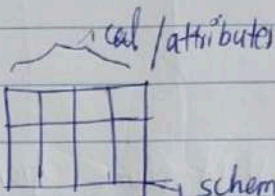 ├→ Unstructured Data → Blob — ADLS
 └→ Semi-structured Data → or Cosmos DB
                                        — No file system
                                        — Flat file

Schema on write          │  Schema on read    — non-relational
                                                      data.
Relational    → RDBMS    │  → No.SQL —   Not only SQL
Data                         can store SQL data and also
     → 1) Create the schema   → flexibility  in the schema
        before / load / store data


        Structured data ——→  Semi-structed data.
        ①  ▦

Create table students (Sid number, name varchar(20) ....); Mango.DB
                          → NoSQL — (JSON  Java script object
                                                      Notation )

| id | Name | H₁ | H₂ | H₃ | .... | H₁₀ |
|----|------|----|----|----|------|-----|
| 1  | Anu  | -  | -  | 2  | ?    |     |
| 2  | -    | -  | -  | ?  | ?    | -   |
| 3  | -    | -  | -  | ?  | ?    |     |

→ only have H₃ all are empty so waste of storage

Limitation are in the back

1) Every row/record in a table have same no. of attribute

→ 10 attributes

→ 10 user_details

| U.id | Name | Age | Profession | | |
|------|------|-----|------------|--|--|
| | | | | | |

Select U-Id, Name from uecdetails
   WHERE age >= 20;

it process in the Traditinal way (because it check all attributes
                                    it was not necessary).

② SQL is not an efficient for reading

NoSQL    Uid [                    ]    ↱ it change column to
Mango-DB  Name [                    ]      rows
→ JSON    Age [                    ]    ↳ Columar Data store
Java script object
Notation

{
  U-id :
  Name :
  hobby : [ " ", " ", " ]
}

{
  U-id :
  Nam :
  hobby : [ ......... ];
}
  ↳ wecan take how many
     you want

BigData
      ↳ processing . ADLS → Big Data processing
required to have
    file system

## Flat Files
=

Blob

&ast; we need not to enable

Storage Account → ☒ to enable
                    to hir namespace

Container — space
↑
Store our data
☐

## BigData processing
=

ADLS

→ Storage
☒ to enable the hir
   namespace
↑
Container
↑
Can create a diretory fil
↑
Store our data
⛁

Access
  Hot — Use more
  cool — Not use more
  Archive —

## Storage services

Structured — Azure SQL $\overset{\text{DB}}{\underset{\text{-SQL server}}{}}$

Semi structured — Cosmos DB → mango DB

un-structured → blob & ADLS    JSON doc

## Azure Databricks

Big data → Hands on

Spark → run all spark job — Introduce the
service Name.
→ programing — Scala
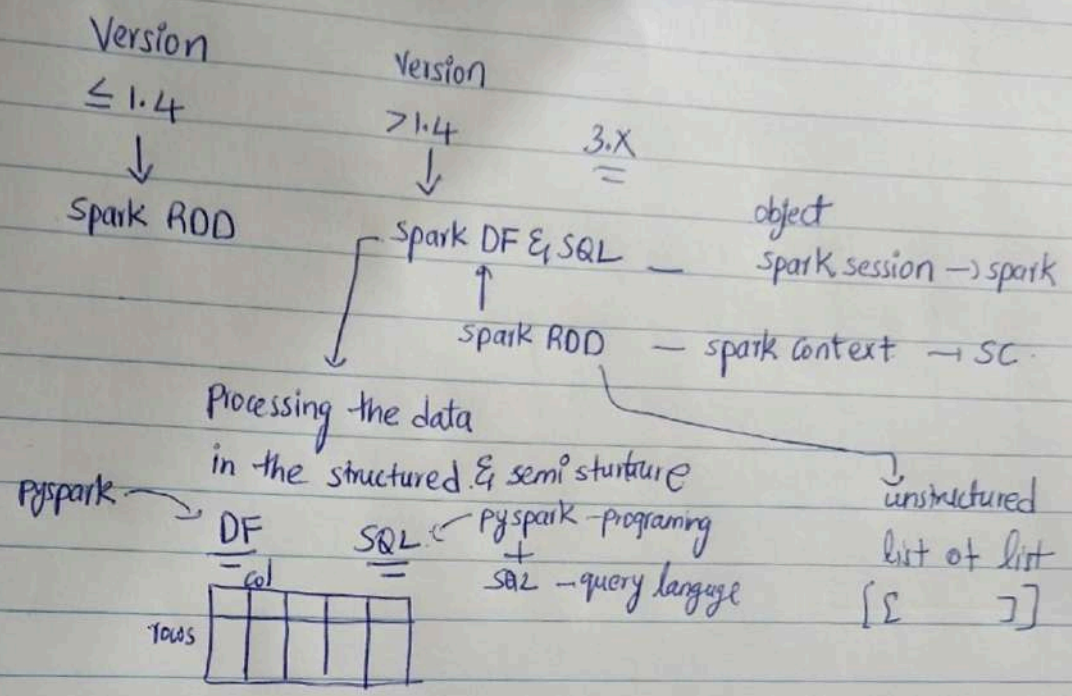Language — Pyspark         Azure Databricks

→ platform - spark - 3. X-version

→ service - delta. bricks

## Azure Databricks

1) Create a Databricks workspace

2) Create a cluster — single Node Cluster

3) Create a Notebook — write spyspark code

4) DBFS → Data bricks file system
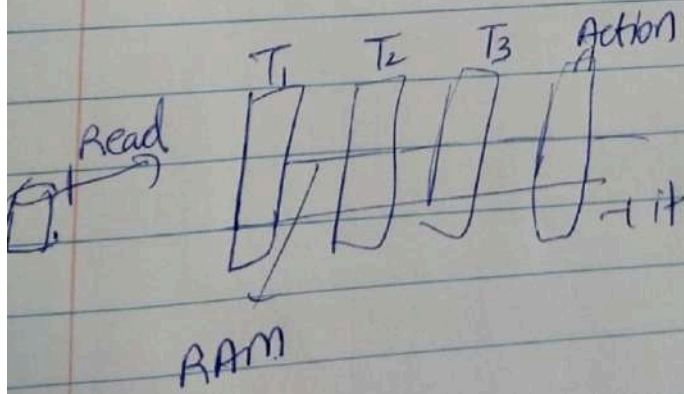    L create DBFS & storege datasets

Spark

Version $\leq 1.4$ → Spark ROD

Version $> 1.4$ → Spark DF & SQL — object  
Spark session → spark

3.X =

Spark ROD — spark context → SC

Processing the data in the structured & semi sturture

unstructured  
list of list  
[ [     ] ]

Pyspark → DF    SQL ← Pyspark -programing  
col    +   SQL → query language  
rows

Spark

2option

not easy memory

1) Tranformation → a, b, c → it will not execute  
↓ Fuction  
[ 'a', 'b', 'c' ]

2) Action

T1   T2   T3   Action

Read

→ it Action Triggerd then only all transform will run

RAM