

INDEX

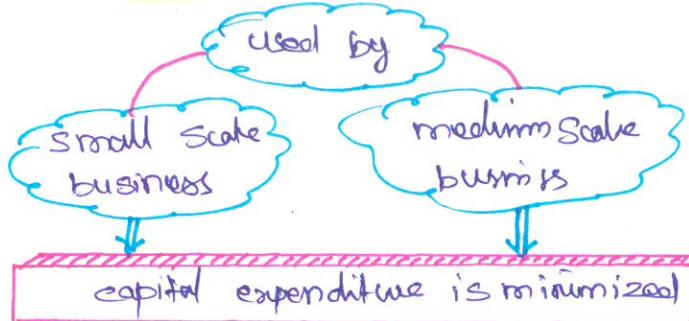
S. NO	Topic	Page No
CLOUD AND SERVICES		
1	Introduction – Evolution of cloud computing <ul style="list-style-type: none"> • Service Oriented Computing • Grid Computing • Distributed Computing • Parallel Computing 	1
2	Hardware Evolution <ul style="list-style-type: none"> • 1st Generation, 2nd Second-Generation , 3rd Generation & 4th Generation Computers 	1
3	Internet Software Evolution <ul style="list-style-type: none"> • Establishing a Common Protocol for the Internet-IPv6 • Building a Common Interface to the Internet 	1
4	Server Virtualization <ul style="list-style-type: none"> • Parallel Processing, Vector Processing • Symmetric Multiprocessing Systems • Massively Parallel Processing Systems 	2
5	Web services overview <ul style="list-style-type: none"> • SOAP ,HTTP & UDDI 	2
6	Infrastructure-as-a-Service (IaaS)	3
7	Platform-as-a-Service (PaaS)	3
8	Software-as-a-Service (SaaS)	3
9	Anything as a service(XaaS)	3
VIRTUALIZATION FUNDAMENTALS		
1	Virtualization architecture	
2	Hypervisors <ul style="list-style-type: none"> • Types of Hypervisors <ul style="list-style-type: none"> ▪ Type-1 & Type-2 • Virtual Machine ,Virtual Network • Virtualization 	3
3	Security Problems in virtualization <ul style="list-style-type: none"> • Data security, Application Security • Virtual Machine Security 	4
4	Cloud Datacenter architecture <ul style="list-style-type: none"> • Data Center Network (DCN) 	4
5	Software defined datacenter <ul style="list-style-type: none"> • XEN Architecture & Microsoft Hyper V 	4
6	Service Oriented architecture <ul style="list-style-type: none"> • Service provider, Service consumer • Service locator 	5
7	Cloud Federation <ul style="list-style-type: none"> • Properties of cloud Federation • Usage of Cloud Federation Properties • Single Sign On(SSO) 	5

8	Presence <ul style="list-style-type: none"> • Presence Protocol 	5
9	Identity Privacy	6
ACCESS TO CLOUD		
1	Hardware and infrastructure	
2	Clients <ul style="list-style-type: none"> • Mobile, Thin & Thick 	7
3	Security <ul style="list-style-type: none"> • Data Leakage, Forensics 	7
4	Network <ul style="list-style-type: none"> • Basic Public Internet, • Providers & consumers 	7
5	Services	
6	Accessing the cloud- Platforms <ul style="list-style-type: none"> • Web application Frameworks • Web Hosting Services, • Proprietary Models 	8
7	Web Applications <ul style="list-style-type: none"> • Sample web applications 	8
8	Web APIs <ul style="list-style-type: none"> • API and API Creators 	8
9	Web browsers	
10	Cloud storage overview <ul style="list-style-type: none"> • The Basics, • Storage as a service • Providers 	8
11	Storage Area Network (SAN) Architecture <ul style="list-style-type: none"> • Types of SAN, SAN Components 	9
12	Collaboration within the cloud <ul style="list-style-type: none"> • Benefits of Cloud Collaboration • Top 5 Features 	9
13	Standards <ul style="list-style-type: none"> • Applications & Clients • Infrastructure & Services 	9
14	Driving forces <ul style="list-style-type: none"> • Popularity, Benefits & Economic Impacts 	10
15	Software and services <ul style="list-style-type: none"> • Vendors & Providers • Mobile Device Integration 	10
16	Developing applications <ul style="list-style-type: none"> • Google & Microsoft 	10
FUNDAMENTALS OF BIG DATA		
1	Introduction	
2	Distributed file system <ul style="list-style-type: none"> • Andrew File System, Google File System 	11
3	Structured and unstructured data <ul style="list-style-type: none"> • Structured data, Unstructured data 	

4	Big data and its importance	
5	5Vs	12
6	Drivers for Big data	
7	Big data analytics <ul style="list-style-type: none"> ➢ Hadoop, Cassandra , Spark 	12
8	Data appliances and Integration tools <ul style="list-style-type: none"> ➢ Benefits of Data Integration ➢ Data Integration Tools 	12
9	Big Data Security <ul style="list-style-type: none"> • Digital Identity Protection • Big Data in Hacking, Steganography 	13
10	Issues in Big data <ul style="list-style-type: none"> • Securing the data, Data Growth Issues • Data Integration 	13
11	Big data applications <ul style="list-style-type: none"> • Applications and Examples 	13
HADOOP AND MAPREDUCE ARCHITECTURE		
1	Big data	
2	Apache Hadoop & Hadoop EcoSystem <ul style="list-style-type: none"> • HDFS , Map-Reduce & YARN 	14
3	Analyzing data with Hadoop streaming <ul style="list-style-type: none"> • Hadoop Streaming • Flow of Data Analysis 	14
4	HDFS concept <ul style="list-style-type: none"> • Introduction & HDFS Architecture • Working of HDFS 	15
5	Interface to HDFS <ul style="list-style-type: none"> • Using Vertica • Using Java 	15
6	Moving Data in and out of Hadoop <ul style="list-style-type: none"> • Hadoop data Ingress and egress • Key elements in data movements 	15
7	Introduction to Map-Reduce <ul style="list-style-type: none"> • Map-Reduce Framework • Word count example 	16
8	Map-Reduce Algorithm and Architecture <ul style="list-style-type: none"> • Architecture & Algorithm 	16
9	Understanding inputs and outputs of Map-Reduce <ul style="list-style-type: none"> • Inputs and Outputs 	17
10	Anatomy of Map-Reduce Job run <ul style="list-style-type: none"> • Mapper, Shuffle & Reducer 	17
11	Failures in classical Map-Reduce and YARN <ul style="list-style-type: none"> • Types of Failures 	18
12	Job scheduling-Data Serialization <ul style="list-style-type: none"> • Types of Schedulers & RPC 	18

Introduction:- to cloud computing :-

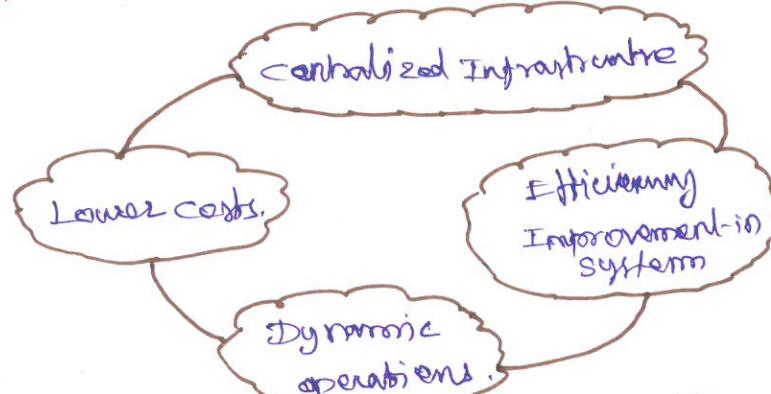
- * Services offering to specific - **Consumers**



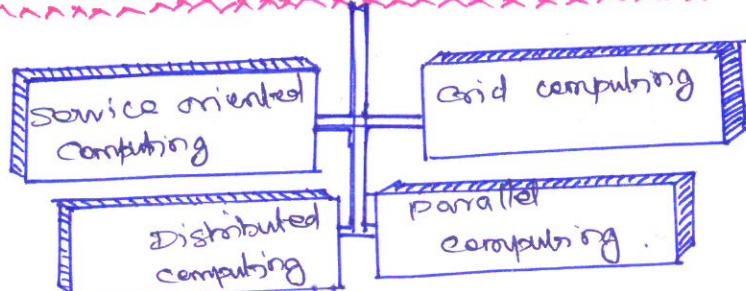
- * Customer access services using



Benefits:-



EVOLUTION OF CLOUD COMPUTING :-



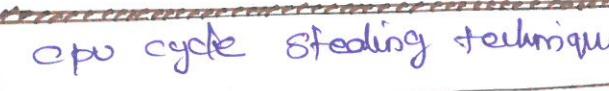
Service oriented computing :-

- * Paradigm for distributed computing
- * built using Services as fundamental

Grid computing :-

- * form of distributed computing

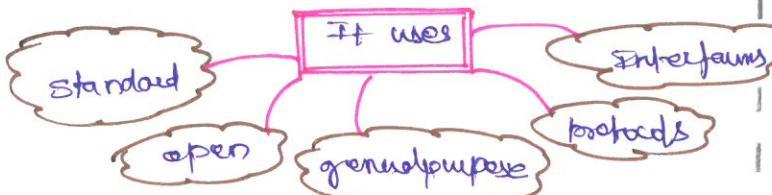
- * Working on the following technique



- * composed of clusters of network,

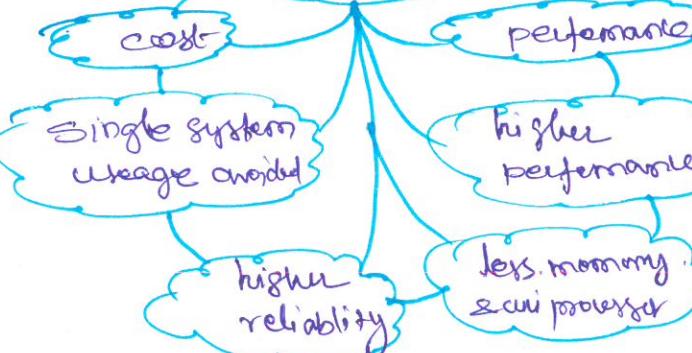
Loosely coupled computers to perform large tasks.

- * Resources are **not** centralized control

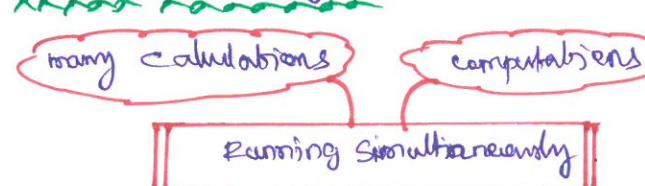


Distributed Computing :-

Advantages:-



Parallel computing:-



HARDWARE EVOLUTION :-

First generation computers:-

- * year 1943

built using
hard wired circuits
vacuum tubes

Data stored using

Paper punchcards

Huge number of punchcard needed

monitoring & searching is difficult

server consolidation

migrating network services and applications from multiple computers to single

change the world toward computing and shape the business world.

growth of mobile technology.

Second generation computers:-

- * year 1946

ENIAC - Electronic Numerical Integrator & computer
It is a first digital computer
transistors & printer currents

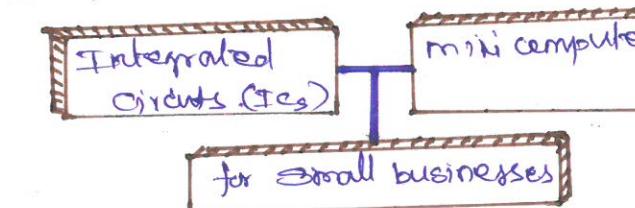
18000 Thermionic valves
Weight 60,000 pounds
25kWph - power consumption

INTERNET SOFTWARE EVOLUTION

Internet -> Internet protocol
Standard protocol used by every computer of internet

Third generation computers:-

- * year 1963

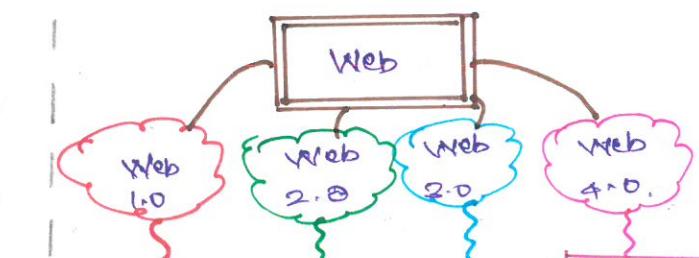


1971 - Intel released first microprocessor - Intel 4004

Fourth generation computers:-

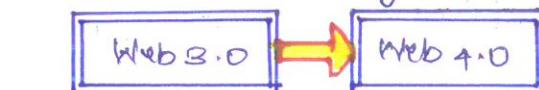
- * year 1974 Introduced a microprocessor with RAM

6000 instructions per second
large scale integration
Very large scale integration LSI & VLSI
performs the functions of CPU



static websites and control
dynamic content in webpage
people centric web, participation, collaboration & non-semantic web
semantic web, that allow user to participate & share the content
Web of goals, goal oriented intelligent resources

- * We are now migrating from



Establishing Common Protocols for Internet:

- * Lower layer protocol layers are provided by Interface Message Processor (IMP) host interface.
- * NCP provided by transport layer.
- ARPANET Host to Host Protocol
- Initial Connection Protocol

- * File transfers : FTP
- * For Email Services : SMTP
- * ARPANET changed from NCP to TCP/IP protocol suite.
- * TCP/IP - Start of Internet.
- * Growth of Internet in 1990s, / Reduction in IP address (IPv4)
- * IETF released IPv6 in 1995

Building Common Interface for the Internet:

- * Mosaic - first web browser - available to general public
- * Mosaic - support graphics, sound, and video clips.
- * 1994 - Netscape released its web browser.
- * 1994 - Mozilla 1.0 released
- * 1995 - Microsoft released Internet Explorer as graphic web browser.

SERVER VIRTUALIZATION

- * Virtualization - running multiple independent virtual OS on a single physical computer.

Parallel Processing:

- * Performed by simultaneous execution of program instructions allocated across multiple processors.

* Running a program in less time.

* Is multiprogramming.

Vector Processing:

- * Operating in multitasking manner.
- * allows single instruction to manipulate two arrays of numbers.
- * applications with lens well-formed data, Vector processing was less valuable.

Symmetric Multiprocessing Systems (SMP)

- * address the problem of resource management in master/slave models.

- * same as single processor, multiprogramming platforms.

Massively Parallel Processing Systems (MPP)

- * System with many independent arithmetic units or entire microprocessor
- runs in parallel massive numbers of computers used to solve single problem.

- * Virtual Server can run its own operating system.

- * Virtual Machine (VM) is a digital version of physical computer.

Types of Virtualization:

- Full Virtualization
- Para Virtualization
- OS level Virtualization
- Network Level Virtualization.

Web Services:

- * WSDL - Service & how
- * UDDI - Static → Service Discovery
- * UDDI - Direct → Service Publication.
- * WSDL - Service Description.
- * SOAP - XML based Messaging
- * Network - HTTP, FTP, Email, MQ, IIOP, etc.

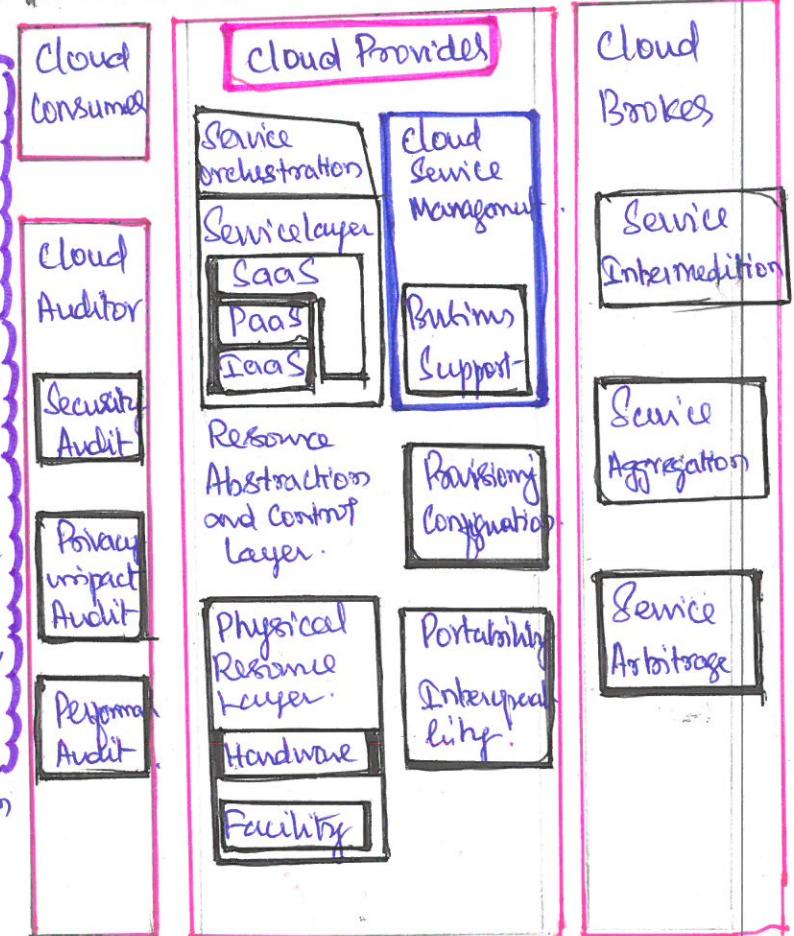
* **SOAP** - Exchange of information in distributed environment.

* **HTTP** - Application layer protocol for transmitting hypermedia documents.

* **UDDI** - Universal Description, Discovery, and Integration is an XML based registry for businesses worldwide to list themselves on the internet.

NIST CLOUD ARCHITECTURE:

- * Cloud computing is a model
 - enable ubiquitous
 - Convenient
 - On-demand network access to shared pool of configurable computing resources.
 - Eg: Networks, Servers, Storage, application, & Services.



Cloud Carriers:

Features of cloud:

- three types of Cloud Service Models.
 - * IaaS (Infrastructure as a Service)
 - * PaaS (Platform as a Service)
 - * SaaS (Software as a Service)

WEB SERVICE OVERVIEW

- ↳ Hardware and software uses HTTP
- ↳ Some protocols respond to request from clients
- ↳ Displays the content of the website

Features of clouds:

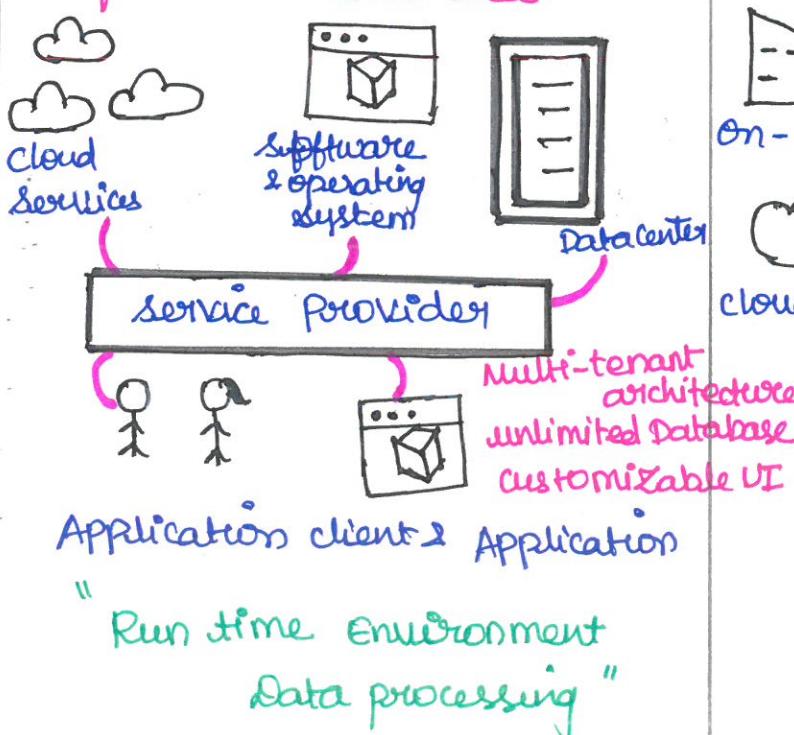
- Infrastructure as a services
- Platform as a services
- Software as a services
- Anything as a services

Infrastructure as a services



Application client (End user) Application

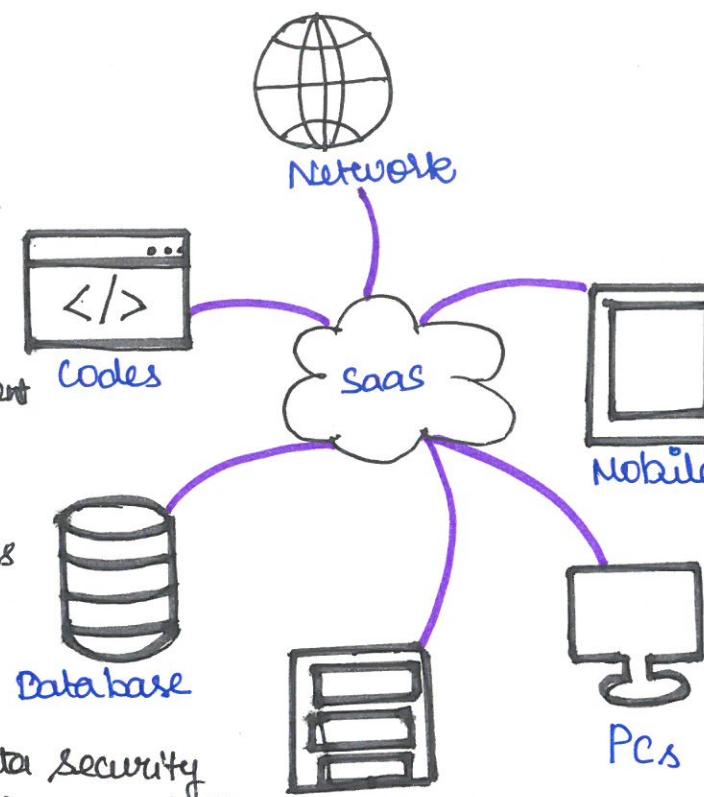
Platform - as - a services



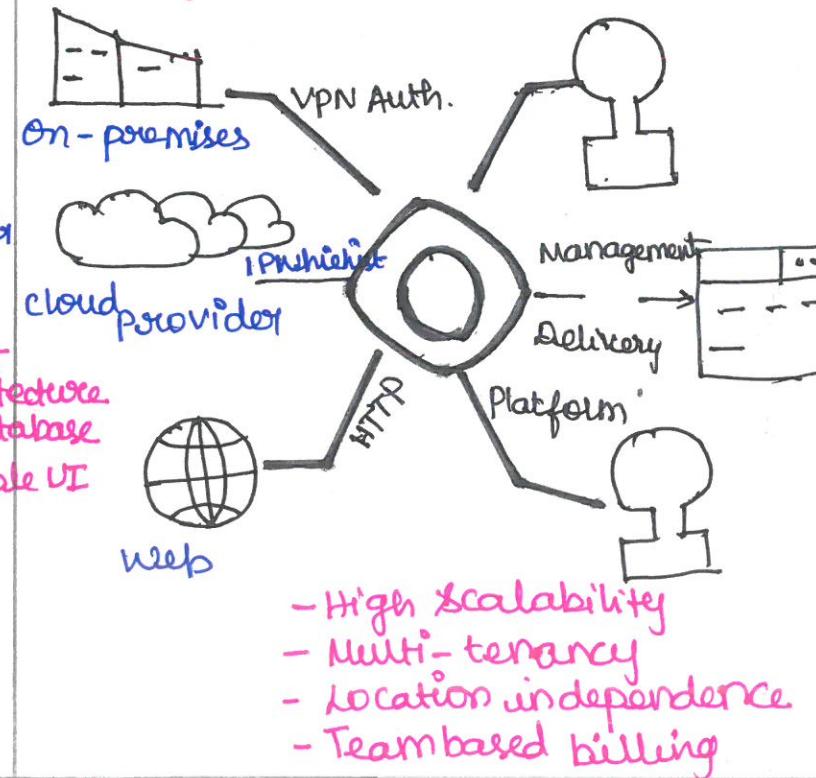
Application client Application
" Run time environment
Data processing "

Software as a service (SaaS)

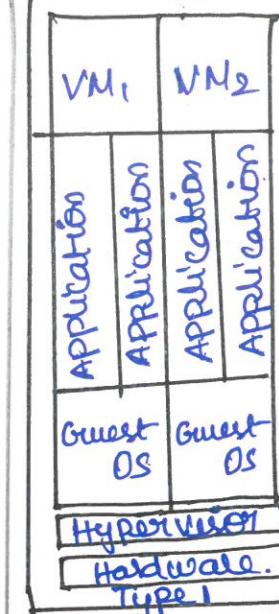
- "scientific application"
- "social network"
- End user application



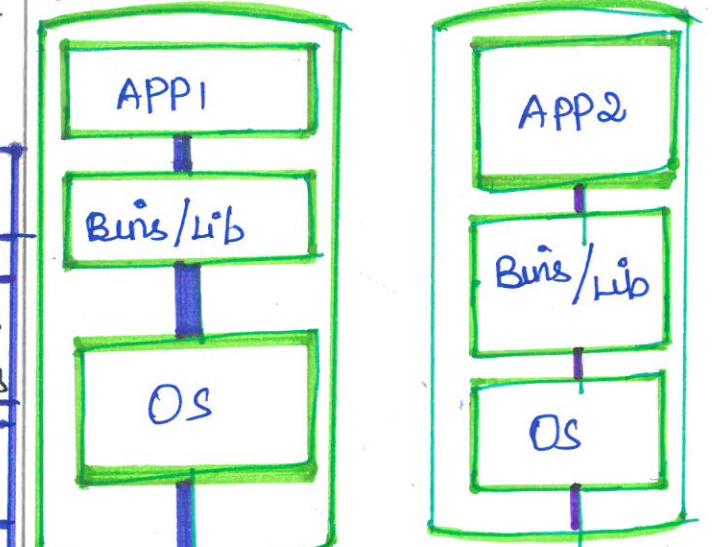
Anything as a services (XaaS)



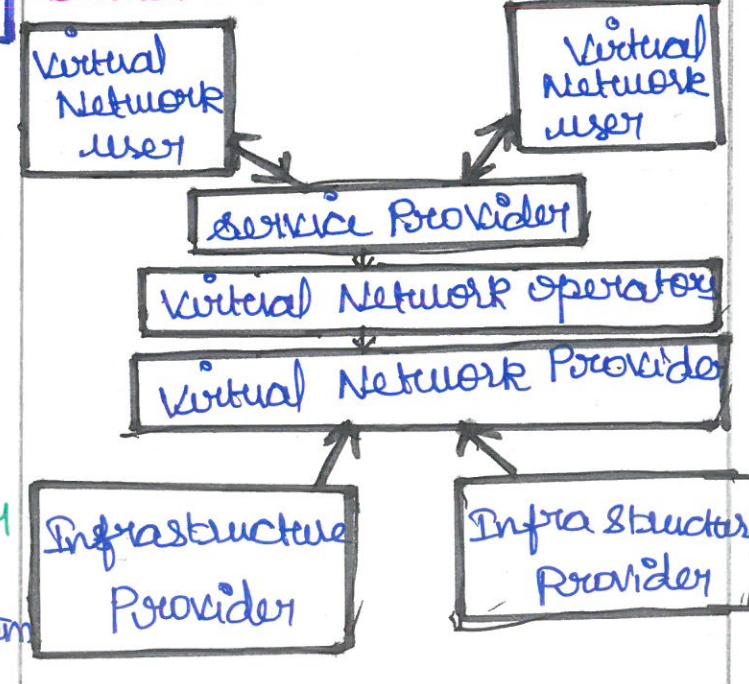
Communicating as a services



Virtual Machine



Virtual Network.

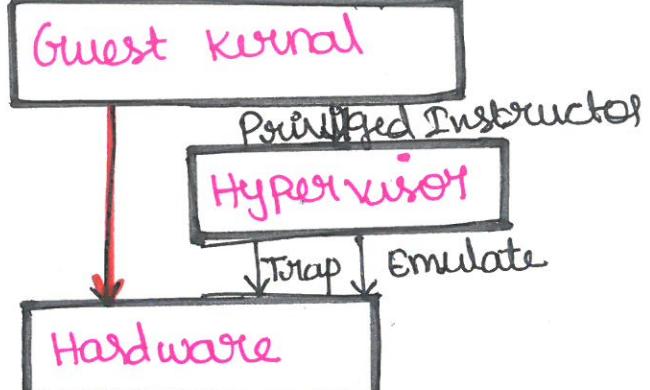


Virtualization

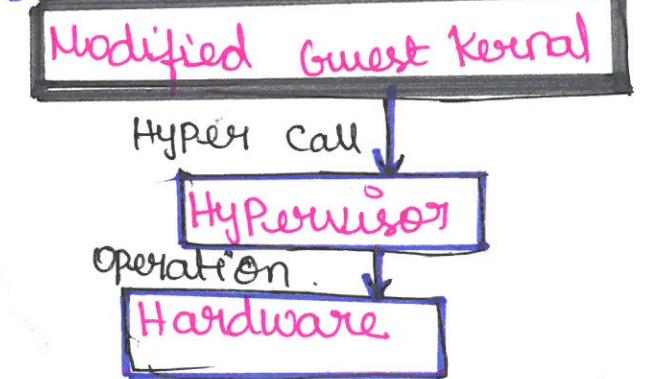
↳ software simulate hardware function and create virtual computer system

1. Full Virtualization
2. Para Virtualization

Full Virtualization



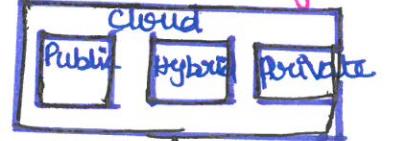
Para Virtualization



Security Problems

↳ insufficient identity, credential, access and key management.

Data Security:



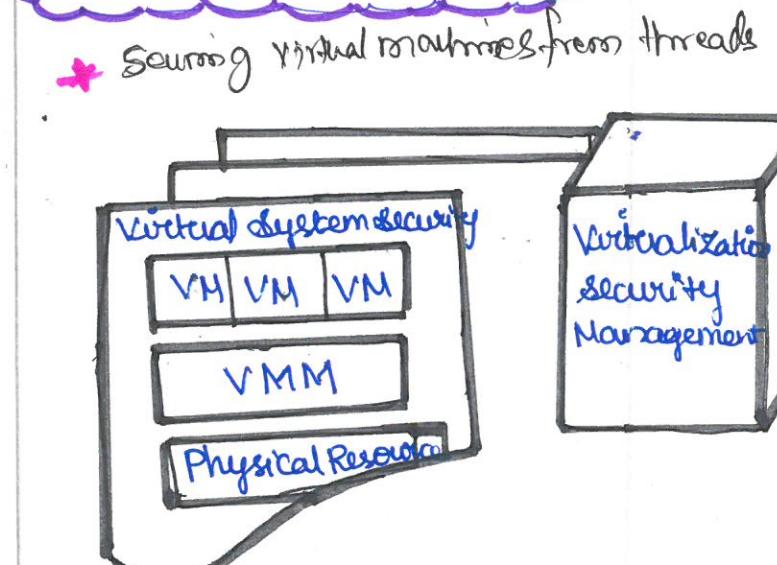
- Data integrity - S/I/W
- Data confidentiality - S/I
- Data availability - H/W
- Data privacy -

Data Security & Privacy

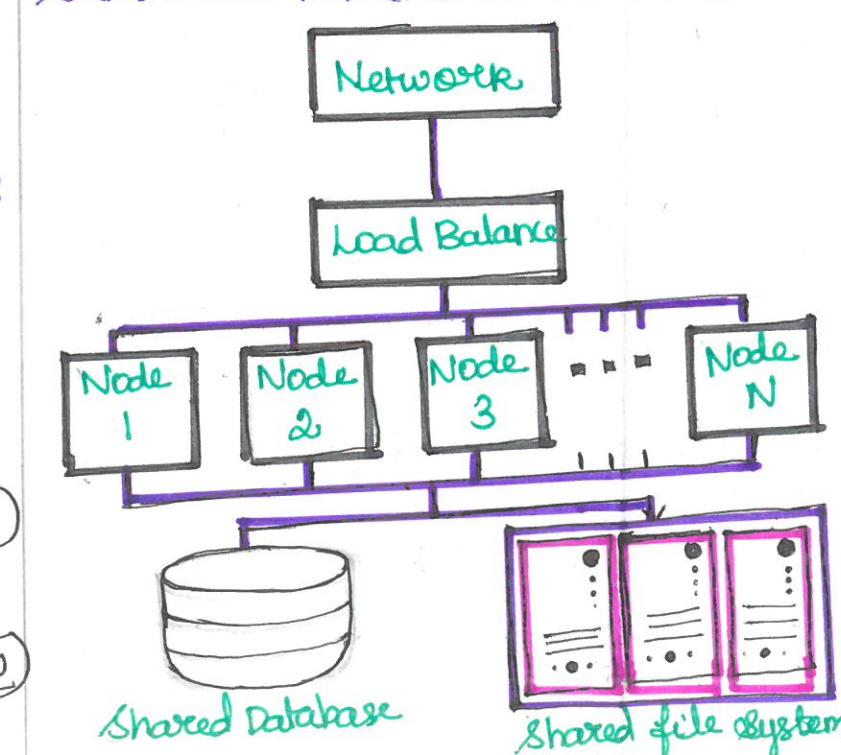
Cloud Application Security



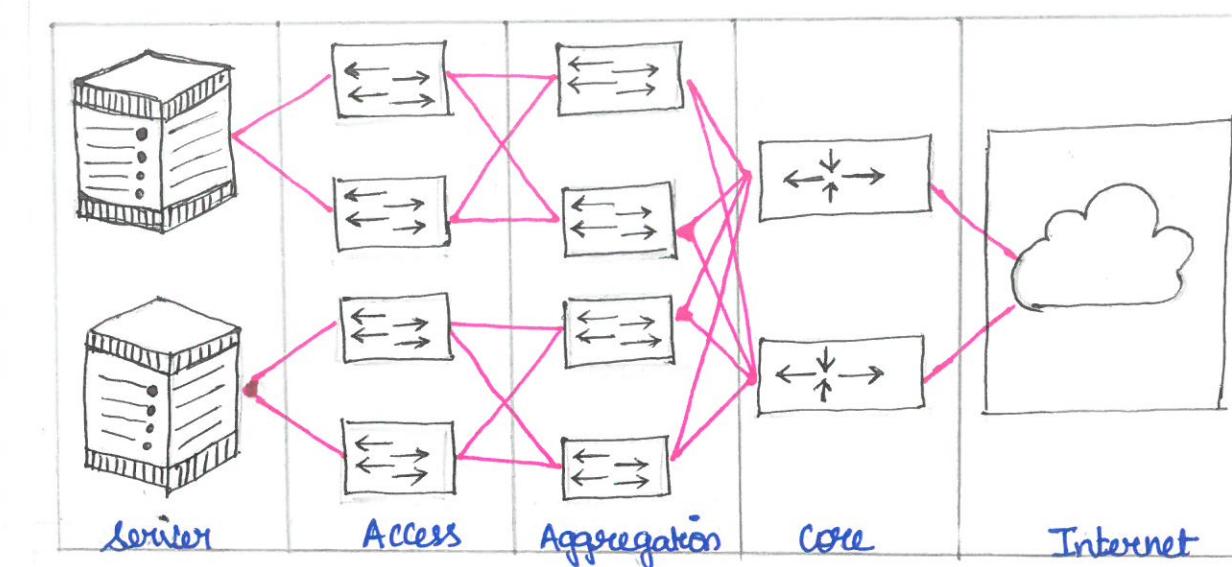
Virtual Machine Security



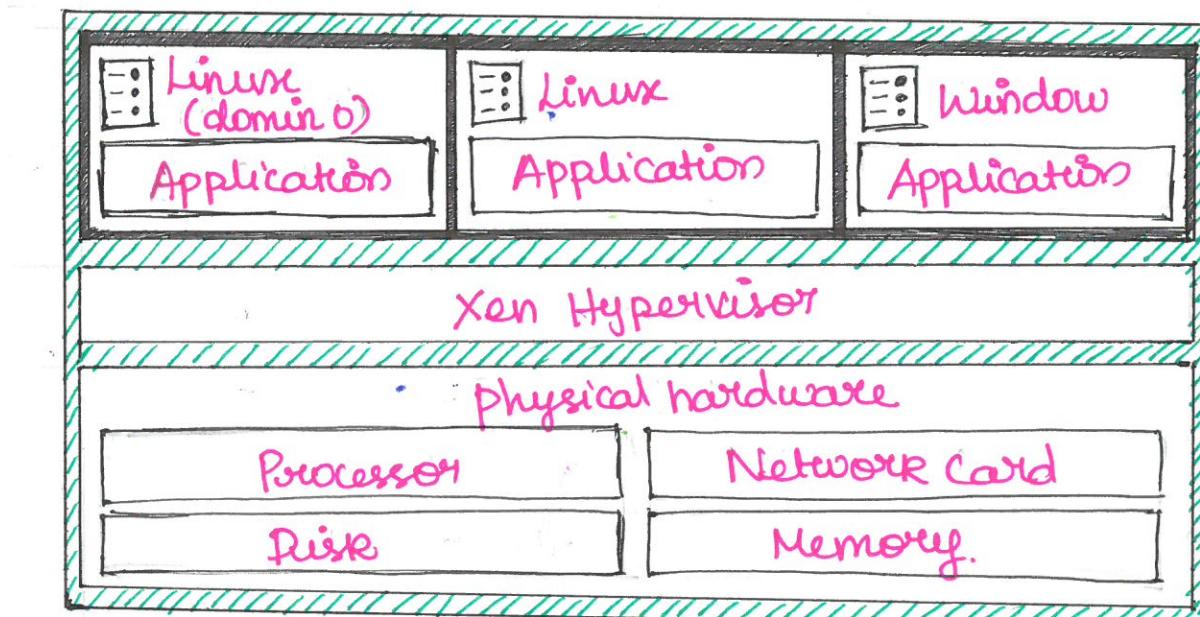
Cloud Data Center Architecture



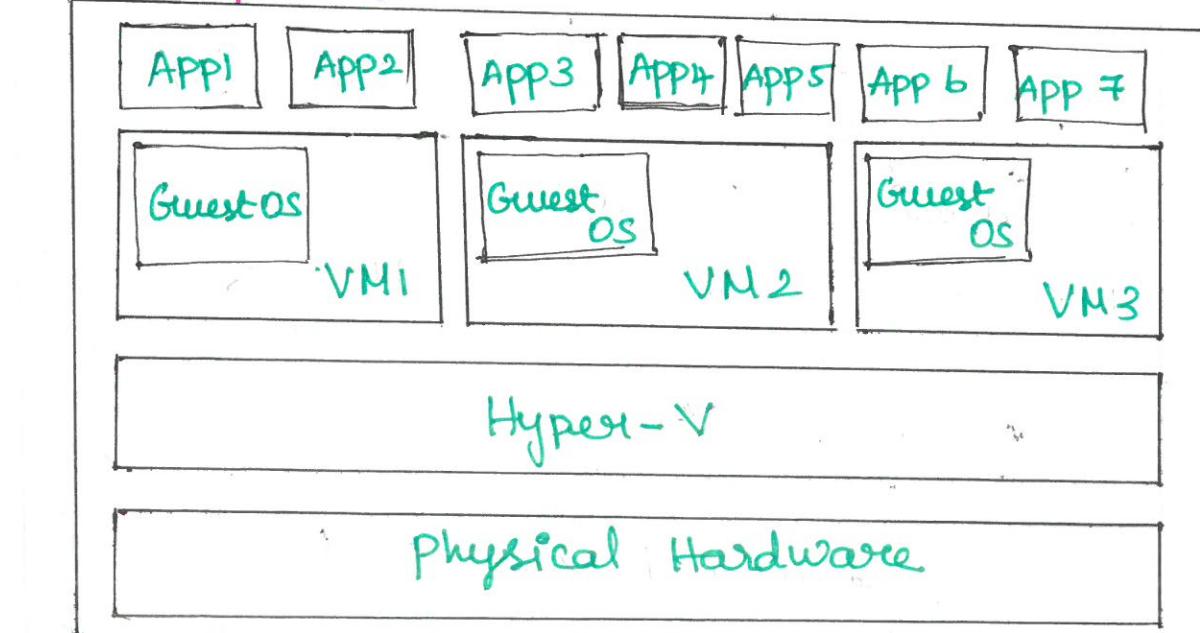
Data Center Network (DCN)



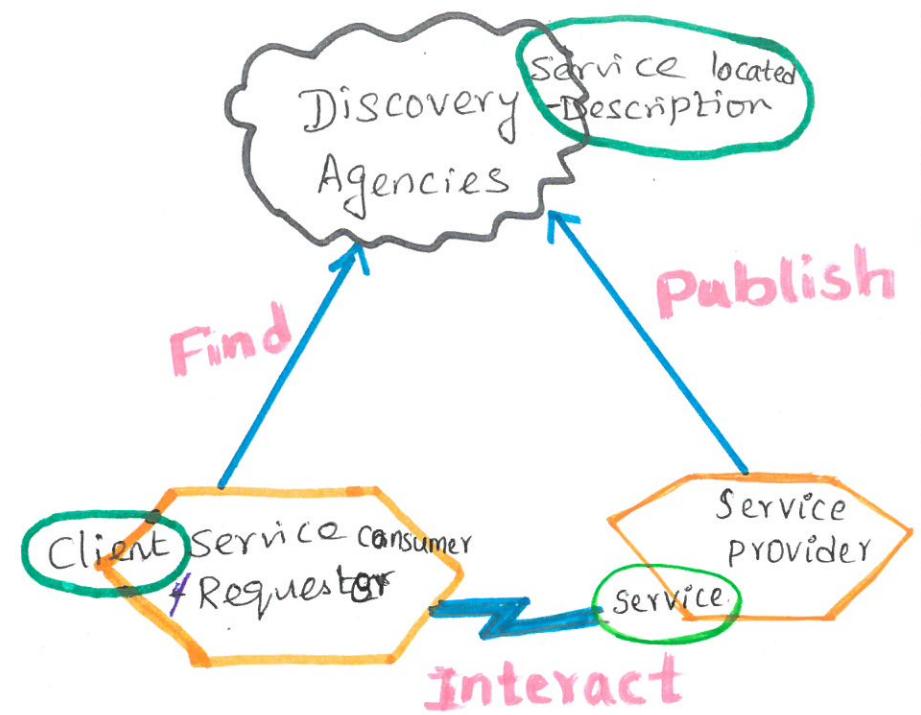
Xen Architecture - Software Defined Datacenter



Microsoft Hyper-V



Service oriented Architecture (SOA) / web service triangle



⇒ Services are the logical entities defined by one or more published interfaces.

Service provider - software entity implements service specification.

Service consumer - calls service provider can be another service (or) end user application.

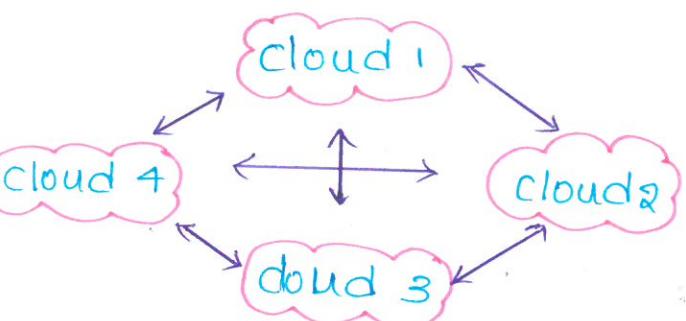
Service locator - Act as a Registry examining Service provider interfaces and Service locations

Quality of Service

Policy, Security, Transaction, Management

Cloud federation :

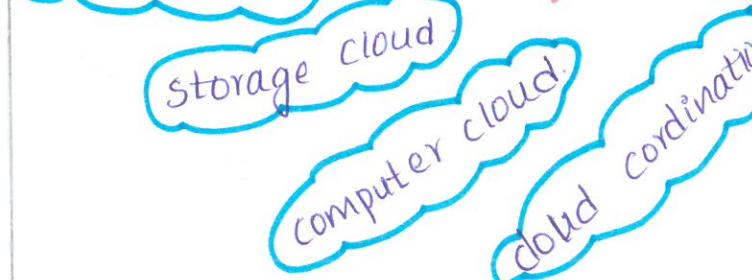
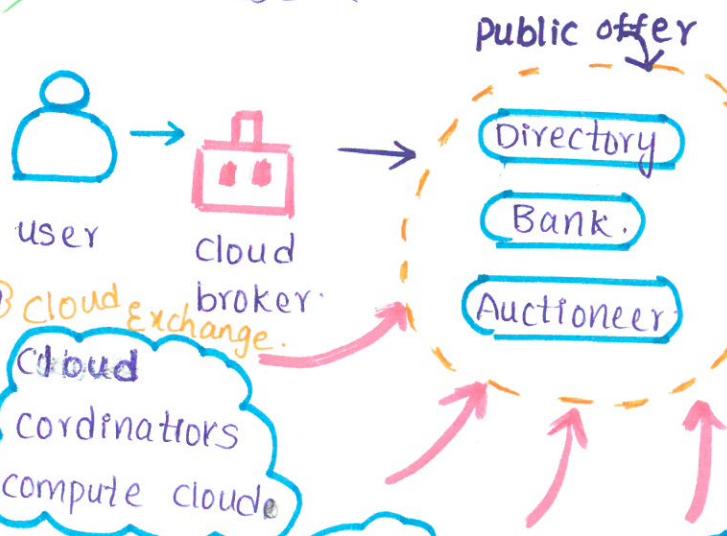
- * Cloud federation also known as federation cloud
- * Deployment and management (External & internal cloud)
- * Multi-National cloud system, integrates private community and public clouds



Cloud federation

⇒ federated cloud is created by connecting the cloud environment of different cloud.

⇒ It can use 2 or more clouds



Properties of cloud federation

- ⇒ users can interact with the architecture either centrally or decentralized manner.
- ⇒ It can be practiced with various niches like commercial and non-commercial infrastructure, software, platform, marketing objects.

Characteristics

- * Authentication
- * Integration
- * Monitoring
- * Object contracts
- * Provisioning
- * Service Management
- * Inter operability
- * Commercialization

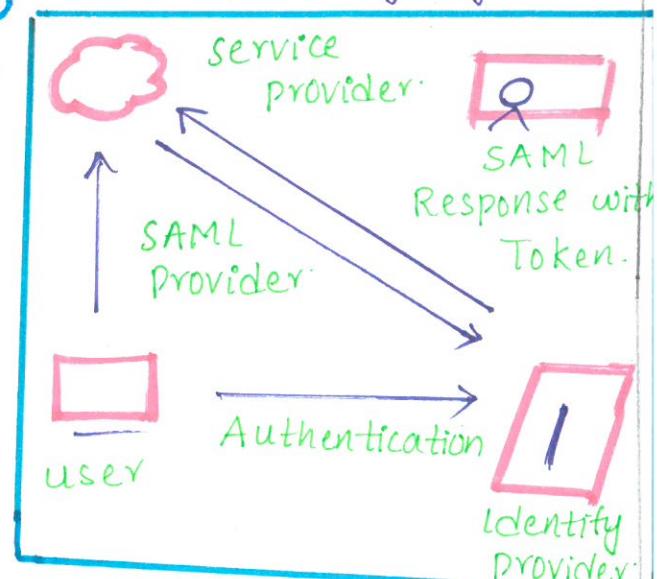
Usage of cloud federation properties

- 1) Interaction of architecture
- 2) Expansion of federation
- 3) Four centric federated cloud
- 4) Business provider → Service customer
- 5) Practice Niche
- 6) Visibility

SSO - single sign on

SSO - It is an authentication method that enables user to securely authenticate with multiple applications.

SAML - Security Assertion Markup Language.



Presence protocol
It works on SSO and steps involved in SSO are as follows :

Step 1 : The service provider sends SAML Request.

Step 2 : If it is not authenticated it will prompt for the authenticated credentials.

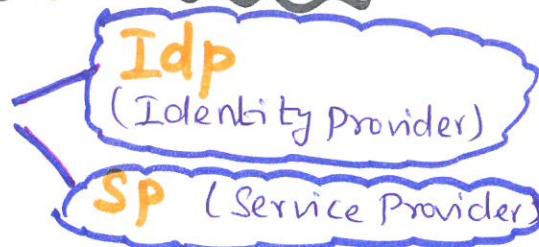
Step 3 : Identify providers, will send an Authentication Token in SAML Response.

Identity

- * Trust between two parties
- * Authenticating users and to convey information needed
- * Authorize accessing resource

Federated Identity Management

TWO Operational roles



Access Management

- * Process of controlling and tracking access
- * Different Privileges within a System on the Individual needs of users

Ex: Payroll System

IAM Service Components [Identity Access Management]

Authentication Services

- * Single sign-on
- * Multi-factor Authentication
- * Session & Token Management

Authorization Services

- * Roles
- * Rules
- * Attributes
- * Privileged access

Directory services

- * Identify Store
- * Directory federation
- * Meta data Synchronization
- * Virtual directory

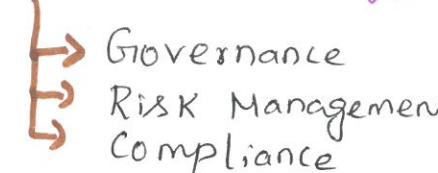
User Management Services

- * Provisioning
- * De-provisioning
- * Self-service
- * Delegation

Identity - as - a - Service (IaaS)

SaaS Model < Solves the identity problem
Provider SSO for web application

Provide Service for Elements



- * IaaS is a prerequisite of Cloud computing for managing the identities and its access.

Languages used: SAML, XACML

Framework: IGF [Identity Governance Framework]

Examples: Ping Identity, Simplified Tricipher & Arcot Systems

PRIVACY

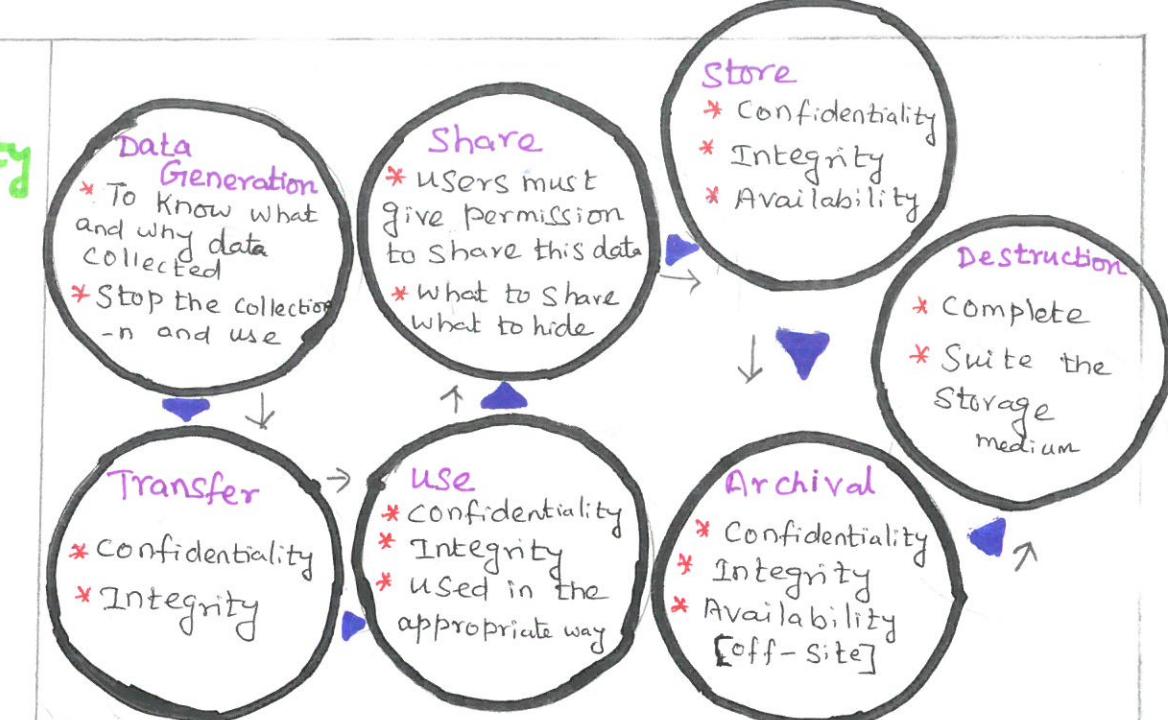
Privacy risk and Cloud

A user's privacy and confidentiality risks vary with terms of service and privacy policy established by cloud provider

- * It is difficult to assess the status of information with legal uncertainties

Protecting - Privacy Information

Federal Trade Commission [FTC] is educating consumers & businesses importance of personal information privacy.



Issues of Privacy Protection

- * Lack of Physical control
- * Protection of Data copies
- * Legal Problems due to privacy laws all over the world

Steps to provide privacy

Attention to privacy of Personal information should be taken in SaaS and managed services

- * Transferring personally identifiable information to and from a consumer's system.
- * Storing personal information on consumer's system
- * Transferring anonymous data from consumer's system
- * Installing software on a consumer's system
- * Storing & processing user data at company
- * Deploying servers.

The Future of privacy in cloud

- * Include better policies and practices by cloud providers

- * Establish Standards

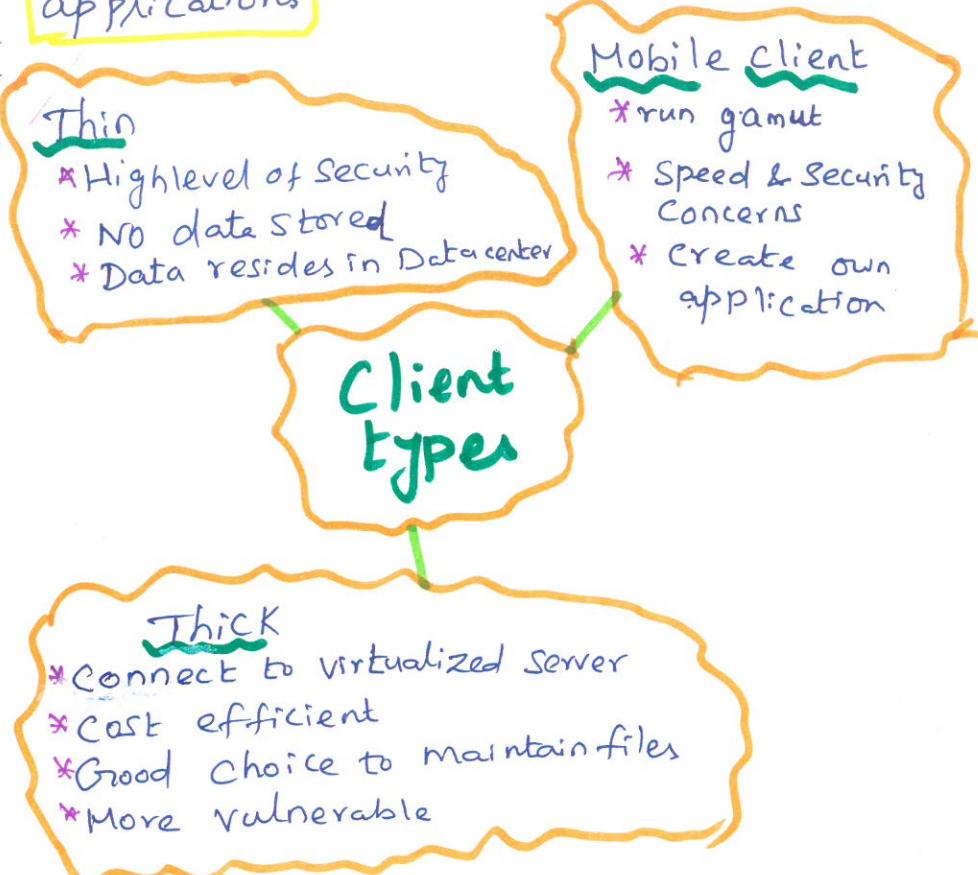
To analyze difference between cloud providers
Assess the risks faced by users

- * Cloud providers would benefit users with greater transparency about risks.

Hardware and Infrastructure

Client

different types of clients offer different ways to interact with your data and applications



Security

- * ISSUE in Cloud computing
- * Third Party stores data

Data Leakage

- * Benefit is the centralization
- * Thick Client
 - download files
 - Maintain
- * Thin Client better for data storage

Forensics

Breach, cloud provider respond with less downtime than investigation.

VM cloned for offline analysis

Auditing

- * Periodic examination
- * assess and document vendor's performance

Internal Audits

External Audits

Audit

Internal Revenue Service (IRS) audits

- * Performed by Certified Public Accounting [CPA]

Network

- * Access cloud via Internet

Basic Public Internet

- * Basic choice for cloud

Connectivity

Example: Internet Service Provider, [ISP] via Broad band, Dial-up

Accelerated Internet

- * Advanced application delivery features

Optimized Internet Overlay

Customers to access cloud via public Internet

Site-to-Site VPN

Connect to Service provider directly using private wide Area Network [PWAN]

- * Confidentiality
- * guaranteed bandwidth

Accelerated Internet

Basic Public Internet

Connection Method

Site-to-Site VPN

Optimized Overlay

Cloud Providers

Third party Company Offers

- Cloud based platform
- Infrastructure
- Application
- Storage Services

Example: Amazon web Services [AWS]

- * Performance Improved
- * Bandwidth charge reduced

Cloud Consumers

- * Include Subscribers to cloud services
- users of cloud services

Example: DropBox, Salesforce

file hosting

Storage, Synchronization will performed on remote servers

Services

Depends on cloud provider

Identity

Payments

Mapping

Search

Integration

Identity

- * users
 - * Digital Identity
 - * Single Sign-on Standard
 - * OpenID Authentication
- Example: Google, IBM

Integration

- * One permises infrastructure
- Example Amazon's Simple Queue Service
- * Use Queue Service
- Example BizTalk Services
- * relay Services

Mapping

Services, Google Map, Microsoft virtual Earth provide cloud based functions.

Payments

- * online payments
- * Signup → Accept credit cards → Send money

Search

- * Rich feature
- * Microsoft Live search is used in cloud.
- * Minimal for small organization
- * Searchability is limited to the organizations

ACCESSING THE CLOUD

PLATFORMS

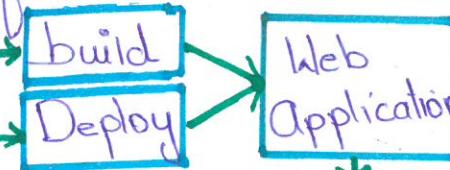
How Cloud Computing Environment is delivered to you

1. WEB APPLICATION FRAMEWORK

Support the development



* Web Application framework Provide a Standard Way to



Types

* Client Side framework

Used to Enhance the User Interface Front End

* SERVER Side framework

Yules You've Set for the Website Back End

Examples: Ruby On Rails, Django, Angular Js, Asp.NET, METEOR

2. WEB HOSTING SERVICE

That Allow you to Store your data and application

Examples:

- * Amazon Elastic Compute Cloud
- * MOSSO

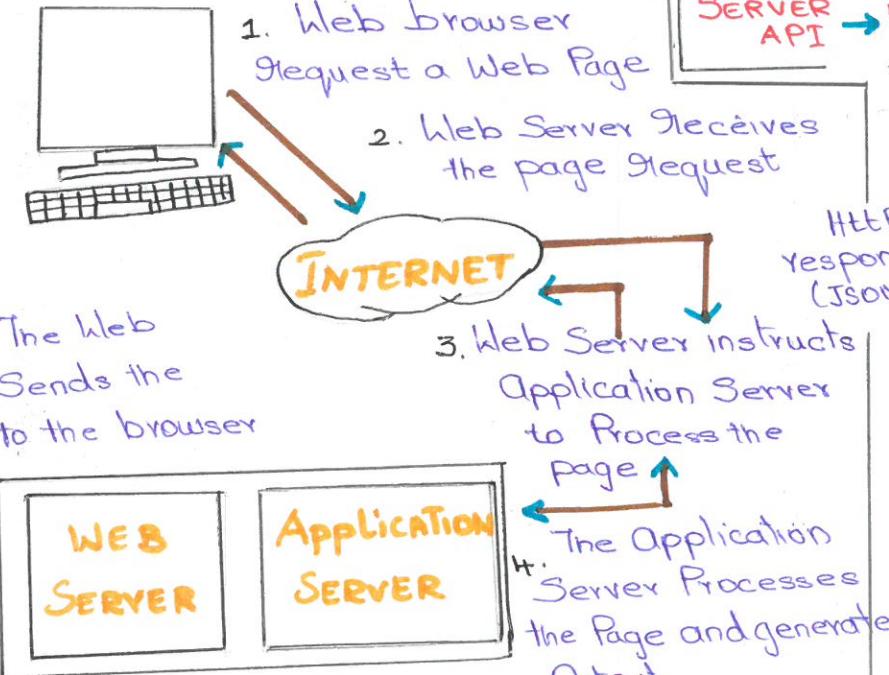
3. PROPRIETARY METHODS

↳ Connecting to the Cloud

Examples Of Own Infrastructure
↳ Microsoft, force.com

WEB APPLICATIONS

* Software or Mobile Application Which Catered Over the Internet And Used in browser



SAMPLE WEB APPLICATIONS

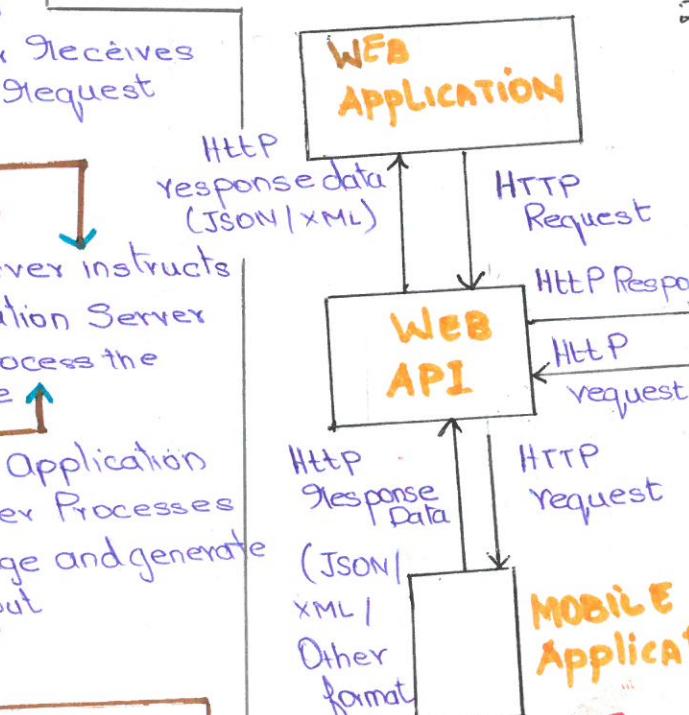
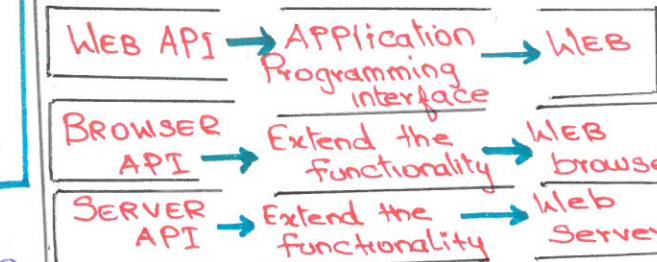
- * Gmail Webmail Services
- * Google Calendar Shared Calendering
- * Google Talk Instant Messaging and Voice OVER IP

Start Page for Creating a Customizable home page On a Specific domain

WEB APIs [Application Programming Interface]

1. What is an API?

- API is a Software Intermediary that allows two applications to talk to Each Other.



[How it Works]

2. API CREATORS

- Google Gadgets
- Google Data API's
- Partnership
- GoGrid
- Apex

1. The Basics

* Some providers are huge And fill an Entire Warehouse

* Niche - Oriented :

- * While, Other store any type Of data
- * Some providers Are Small

WEB BROWSERS

Takes you Anywhere on the internet

1. firefox

- * Smaller, faster, More Secure
- * Protects from -
 - Cross-Site Tracking Cookies
 - fingerprints
 - Cryptominers

2. Chrome

- Open Source Program Accessing the World Wide Web
Running Web-based Applications

3. Internet Explorer

- free graphical browser
- Maintained by Microsoft for legacy Enterprise Uses.

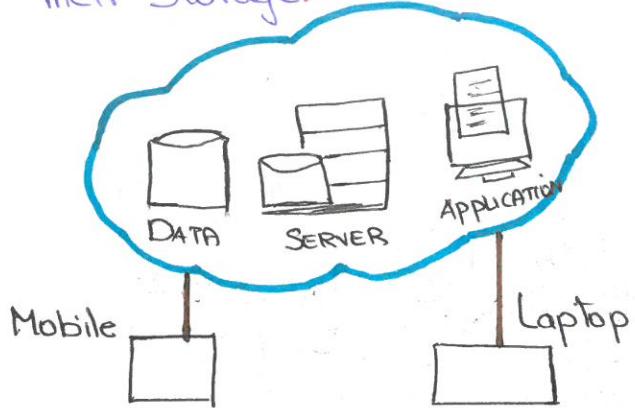
CLOUD STORAGE

Store your Data with Internet Access get it from Any Location.

Store Just Email OR digital Pictures

2. STORAGE As a Service

- * Another name "Software as a Service"
- * Third party Provider Stents Space for their Storage.



ADVANTAGE	DISADVANTAGE
Cheaper upfront	Difficulty Chaining Provider
Easy to implement	Data Stored offsite
Integrated With your Stack	Dependence on the Service infrastructure
minimal Commitment	Compromise Security

3. Provider As a Service

- * Third Party Provider delivers hardware And Software tools Over internet



Encryption and authentication are two Security Measures you can used to keep your data Safe on a cloud storage provider

4. Advantages And Cautions

1. Cost	2. Availability
3. Scalability	
4. Productivity	
5. Upgradeability	
6. Platform Support	

1. Security	2. Control
3. Reliability	
4. Lock in features	
5. Integration	
6. Compatibility	

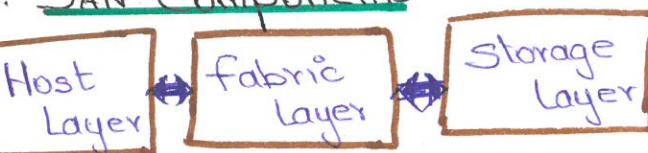
STORAGE AREA NETWORK (SAN)

- Network Of Storage device
- Accessed by multiple Server
- Provide a shared pool Of Storage Space.

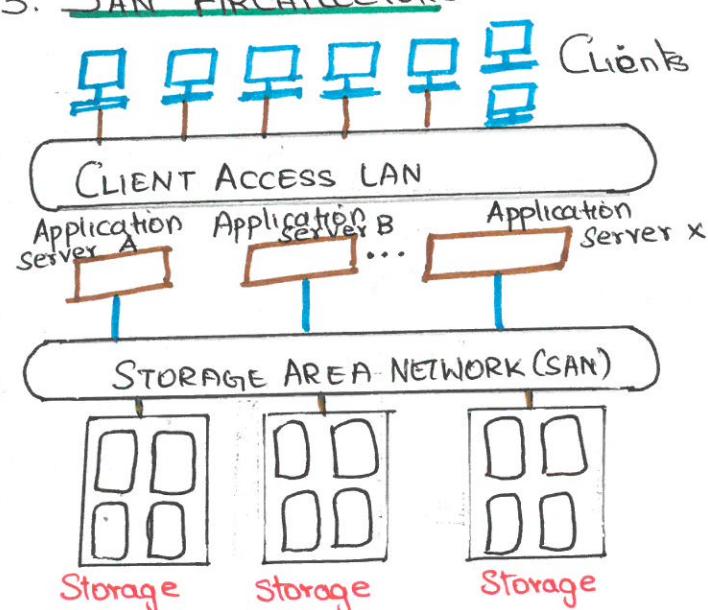
1. Types Of SAN

- fiber Channel Protocol (FCP)
- Internet Small Computer System Interface (iSCSI)
- Fiber Channel Over Ethernet (FCoE)
- Non-Volatile Memory Express Over Fiber Channel (FC-NVMe)

2. SAN Components



3. SAN Architecture



COLLABORATION Within the Cloud

1. Benefits Of Cloud Collaboration

- Reduced investment, Scalability
- Improved Organization
- Higher levels Of Participations
- Large files Are Easy to Access
- Updates in Real time
- Improved Brainstorm

2. Top 5 Cloud Collaboration Providers/Features

- Access files Anytime through the Internet
- Real time Communication
- Setting Custom Permission levels
- Version Control
- Centralized file Storage

STANDARDS

1. Application

- The Protocols that Are Used to Manage Connection between both Parties.

A. Communication:

- Computer's need a Common Way to Speak With One Another

(i) HTTP (Hypertext Transfer Protocol)

- Stateless protocol
- Transfer data between the Cloud and your Organisation

(ii) HTTP 1.1

- Eight Methods to describe how the desired Action to be Performed On the Server.

REQUEST	DESCRIPTION
HEAD	Ask for the Response
GET	Request information from Server
POST	Submits data to be Processed to the Server
PUT	Upload the Resource

Delete

Deletes the Specified Resource

TRACE

Adding or Changing the Request.

Options

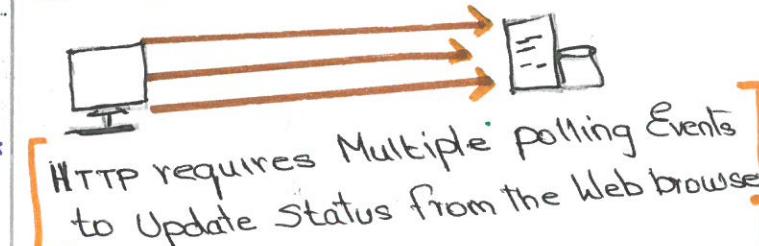
Returns HTTP Methods that the Server Supports for the given URL

Connect

Convert Request Connection to transparent Tcp/IP tunnel.

(iii) XMPP (Extensible Messaging and Presence Protocol)

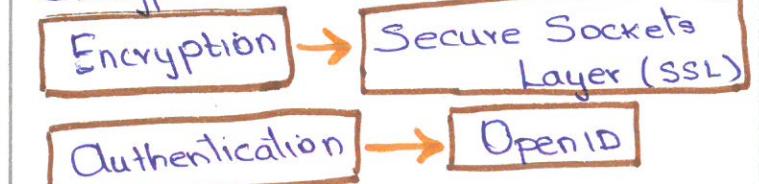
- Allows for two Way Communication And Eliminates Polling.



XMPP maintains a Connection between the Client and the Web Server

B. Security:

- Securing your Cloud Sessions Can be Accomplished Via Encryption And Authentication



(i) SSL: Encryption Link between Web Server and Browser.



(ii) OpenID:

- OpenID is free
- lower Cost for Password and Account Management
- Used to identify Web sites

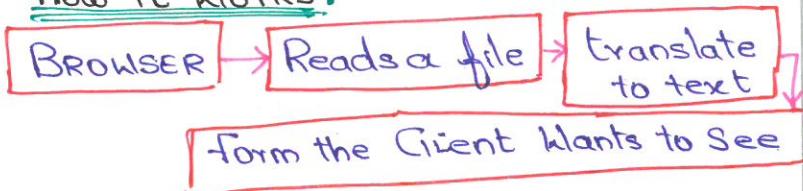
2. Clients:

- It means to store and display information

iii) HTML:

Code in the form of a hyperlink which takes you to another page.

How it works:



iv) Dynamic HTML:

- * Controlling the Standard HTML Codes

There are four parts to DHTML:

- * Document Object Model (DOM)
- * Scripts , * HTML
- * Cascading Style Sheets (CSS)

v) JavaScript

* "Client-Side Web development"

Example:

Opening or popping up new windows and having control of the size and attributes of the window.

3. Infrastructure:

* The Internet: (machine running on a remote server and displayed at your organisation)

* Locally: (having your "clients" session runs on a local server and displayed at their desktops)

vi) Virtualization:



(Application runs on a "Server" and are "displayed" on the Client)
(The Server can be Local or on the Other Side Of the Cloud)

Benefits End Users by:

- * Expanding Virtualization Solutions
- * Expanding interoperability and supportability
- * Accelerated availability of new Virtualization aware technologies

4. Service:

- * Supports Machine to Machine interaction over a network

i) DATA: Stored and served up with number of mechanism.

- JSON (JavaScript Object Notation)
- * Light weight Computer data interchange format
- JSON Basic (subset of Javascript)

Driving forces

1. Popularity:

- * Software Vendors love it
- * Enterprises love it

2. Virtualization Benefits:

- * Easier for independent software vendors to adopt

3. Economic Impact Benefits:

- * Subscription based payment model of "SaaS" makes it more appealing in these tough times

SOFTWARE AND SERVICES

v) VENDORS:

- * Microsoft
- * Apple
- * Adobe
- * Google
- * Salesforce.com
- * WeatherBug
- * Dicentral's Disintegrator EPI Solution

vi) Providers:

- * Creating your own software plus service deployments

vii) Adobe AIR:

- * Adobe Integrated Runtime (AIR)
- * Cross Operation System Application runtime that allows developers to use

- * HTML/CSS , * AJAX , * Adobe Flash
- * Adobe Flex to Extend (RIA)

viii) Apple iPhone SDK [MOBILE DEVICE INTEGRATION]

- * iPhone SDK providers developers with a rich set of APIs
- * iPhone SDK can download free
- * Third party developers are able to build native applications with rich set of APIs.

Developing Applications

i) Google:

- * Getting Started guide
- * Python library documentation
- * Examples showing Python code for accessing force.com
- * Testing harness for the provided library.

* Wiki FAQ page on developer.force.com with best practices and latest tips and tricks

ii) Microsoft:

- * Tool provided for developers who want to write application
- * Runs partially or entirely in a remote data center.

iii) LIVE SERVICE:

- * Set of building blocks within the Azure Service Platform
- * Used to handle user data and application resource.

iv) Microsoft SQL Service:

- * Cloud as a web based, distributed relational database.

- * Store and retrieve structured, unstructured, semi-structured data.

v) Microsoft .NET Services

- * Tool for developing loosely coupled cloud-based application
- * Access to control to help secure application
- * On-premises environments to the cloud

vi) Microsoft Sharepoint Service and Dynamic CRM Service:

- * Allow developers to collaborate and build strong customer relationships
- * Tools like "Visual Studio", "Sharepoint", "CRM Capabilities"

Fundamentals of BIG DATA

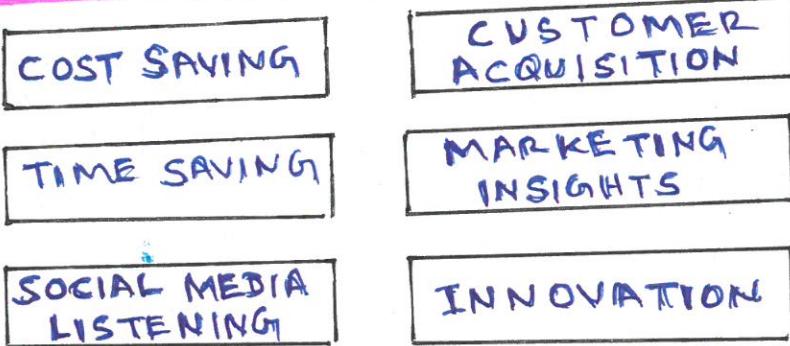
Introduction:-

- Collection of data that is **huge in volume**
- Large size and Complexity**
- Big data** is also a data but with huge size.

Example :-

- New York Stock Exchange
- Social media

Importance of Big data

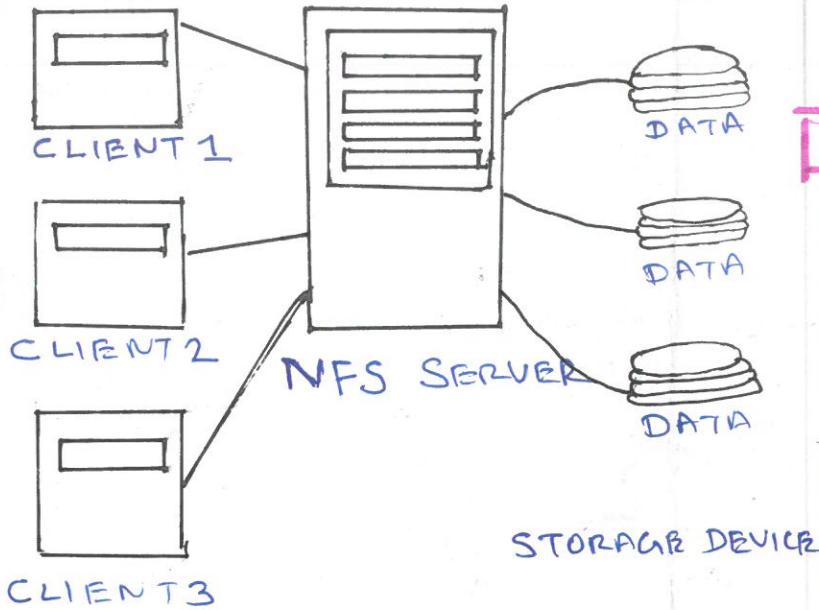


Distributed file System:

- Distributed on various **file servers** and locations
- Access and store **isolated data**
- Access files from any system.
- To share information and files.
- Provides users access control

Two Components of DFS

- Local Transparency
- Redundancy



Features of DFS

- Data sharing of **multiple users**.
- User **mobility**
- Location **transparency**
- Location **independence**
- Backups** and System Monitoring.

Applications of DFS

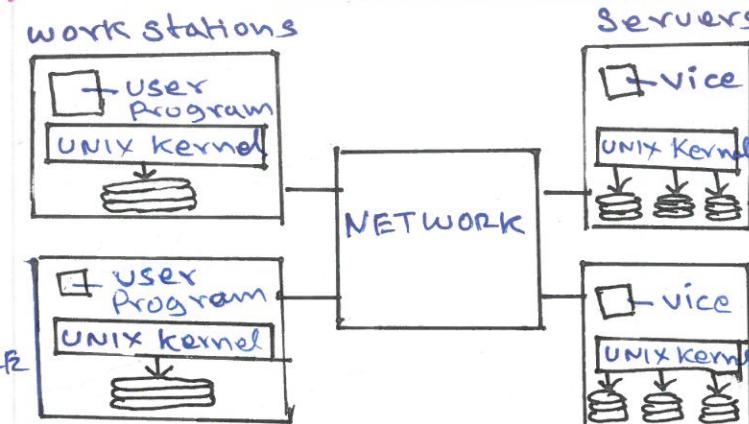
- Hadoop
- NFS (Network File System)
- SMB (Server Message Block)
- Netware
- CIFS (Common Internet file system)

AFS → (Andrew File System)

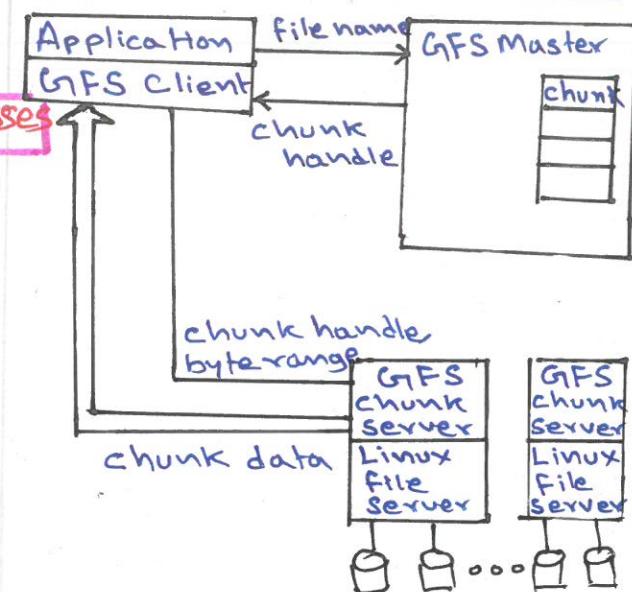
- To support information **sharing** on a large scale
- Transferring whole files between server and client with caching.
- Uniform **name space**
- Location independent file sharing
- Client-side **caching**
- Secure** authentication.

- Replication
- Whole-file serving
- Whole-file Caching

AFS Distribution of Processes



GFS Architecture



Advantages :-

- Allows the user to access and store the data
- To improve the **access time**, network efficiency and availability of files.
- Provides the **transparency** of data.
- Permits the data to be shared **remotely**.
- Helps to change the amount of data and exchange data.

Structured and Unstructured

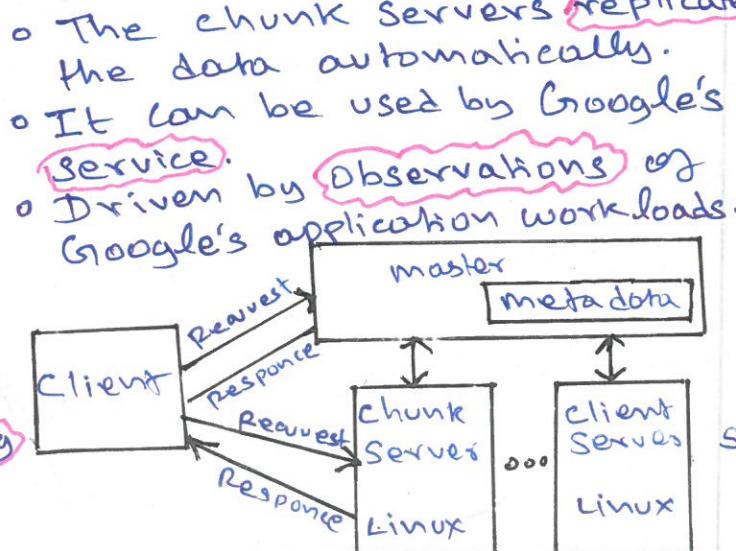
Predefined Formatted data
Examples:-

- Names
- Date
- Address
- Geo Location
- Credit Card Number

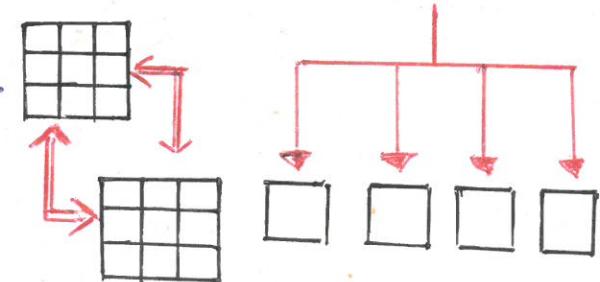
Unstructured data
Example:-

- E-mail
- Social media posts
- IOT sensor data
- Satellite imagery

SQL VS NoSQL



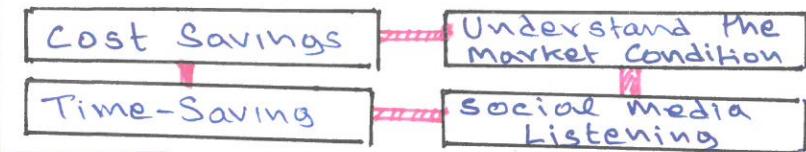
SQL



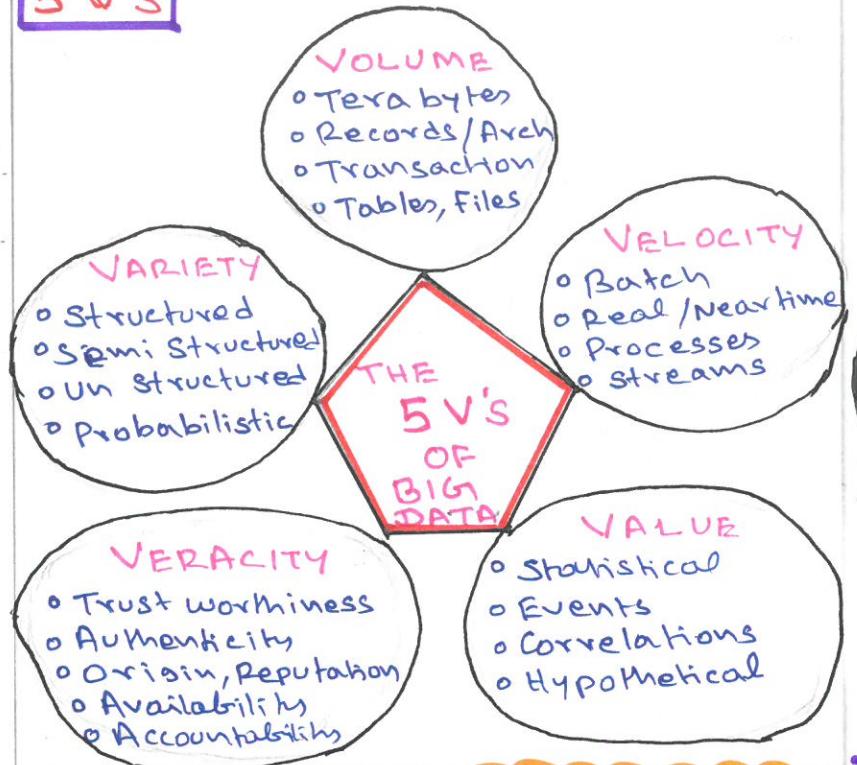
NoSQL

Un-Structured Data

Big data and its importance



5V's



Volume (Huge amount of Data)
Velocity (High speed of accumulation of data)
Variety (Heterogeneous Sources) data
Veracity (Inconsistencies and uncertainty in data)
Value (Worth of data)

Drivers of Big data

Big data has quickly become one of the most desired topic in industry.

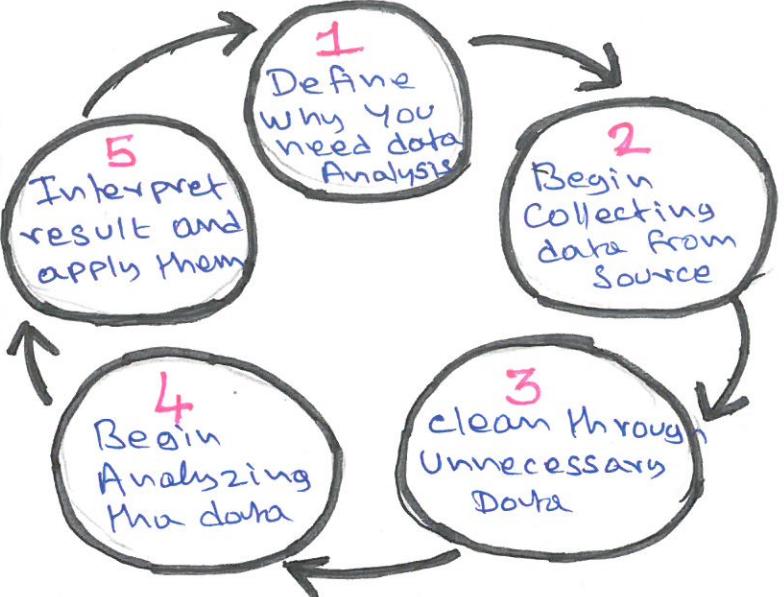
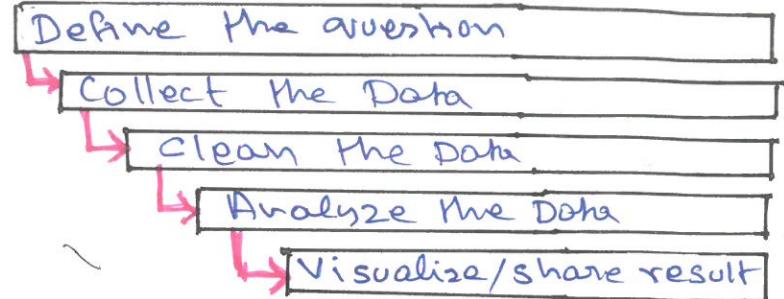
The Digitization of Society

- The Drop in technology costs
- Connectivity through cloud computing
- Increased knowledge for Data Science
- Social media Application
- The rise of Internet of Things

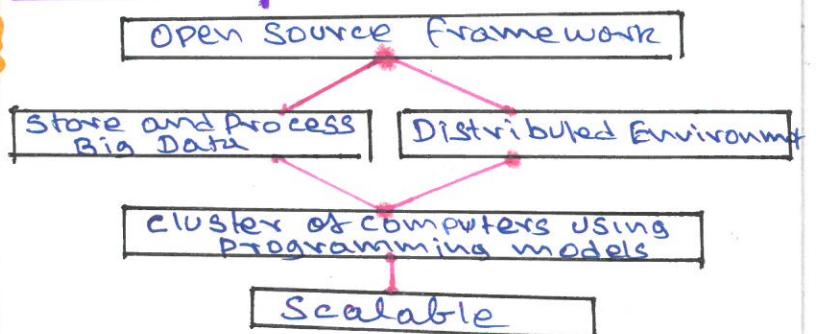
Big data analytics

Introduction
Analysis of large volumes of diverse data sets, using advanced analytic techniques.

Five different steps of the big data analytics Process:



Hadoop



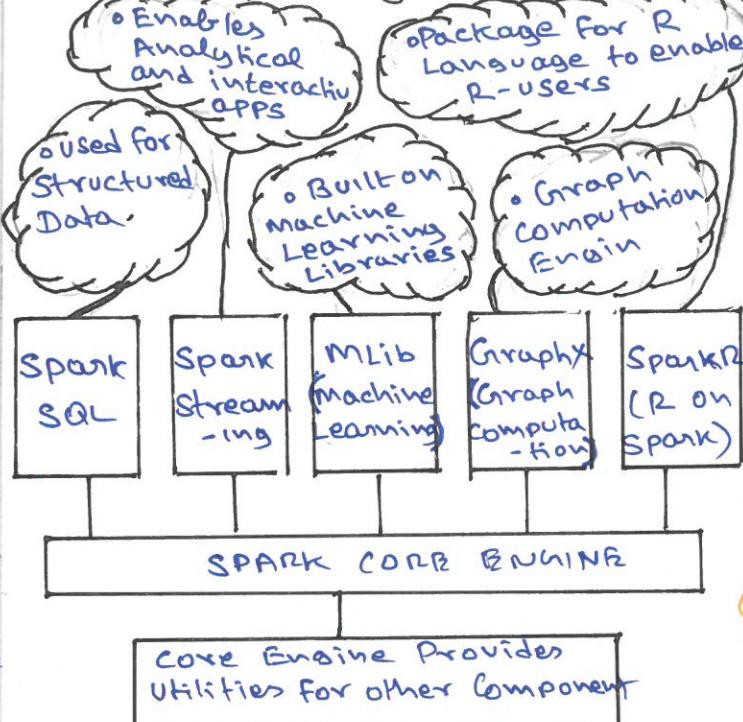
Cassandra

- Open Source
- Distributed and Decentralized
- Distributed Storage System
- Managing large amounts of Data
- Open-Source availability
- Distributed Foot Print
- Scalability.
- Cassandra Query Language
- Fault tolerance
- Schema free

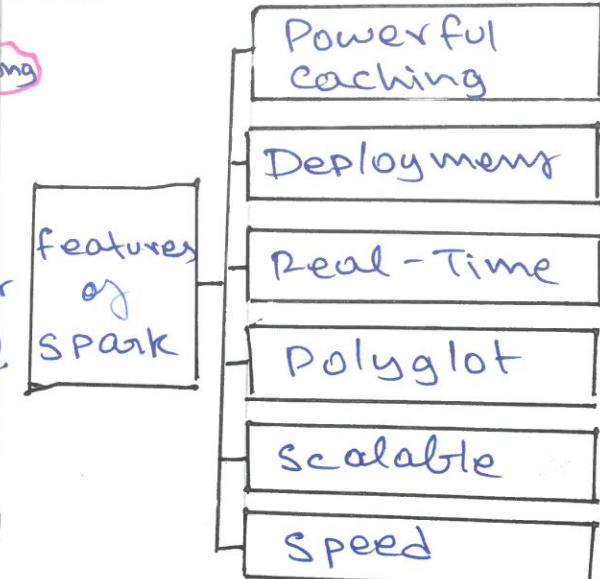
Spark

- Spark is a real time data analytics in distributed Computing
- It executes in memory to increase speed of data Processing.
- Faster for large scale data.
- It requires high processing power
- Resilient Distributed Data(RDD), is data structure of Spark

Components of Spark

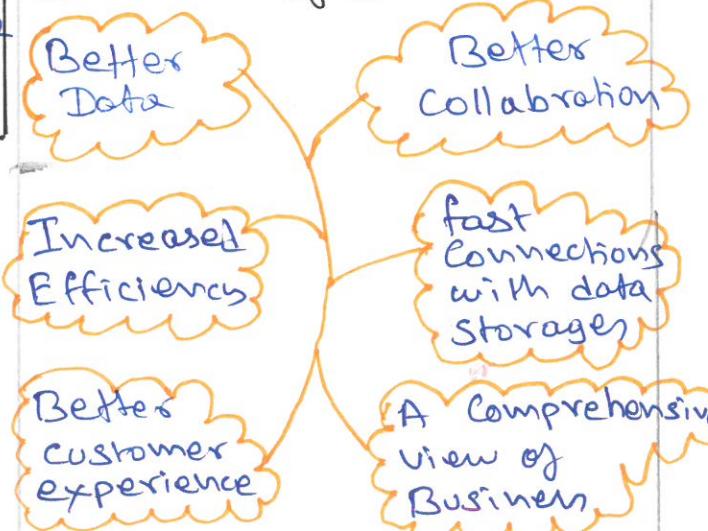


Features of Spark

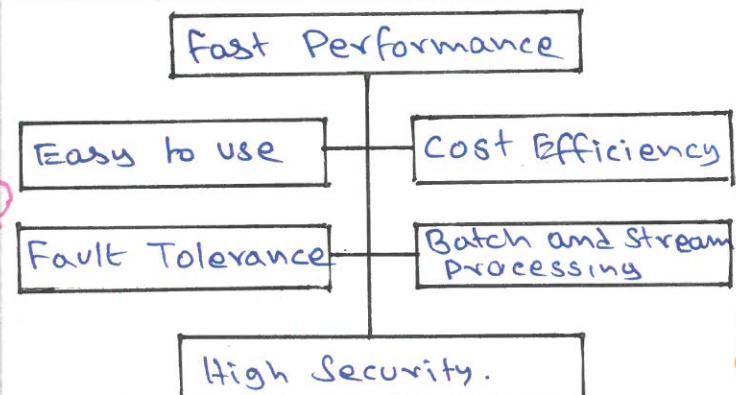


Data Appliances and Integration Tools

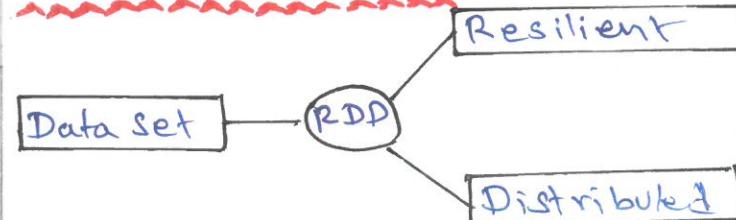
Benefits of data Integration



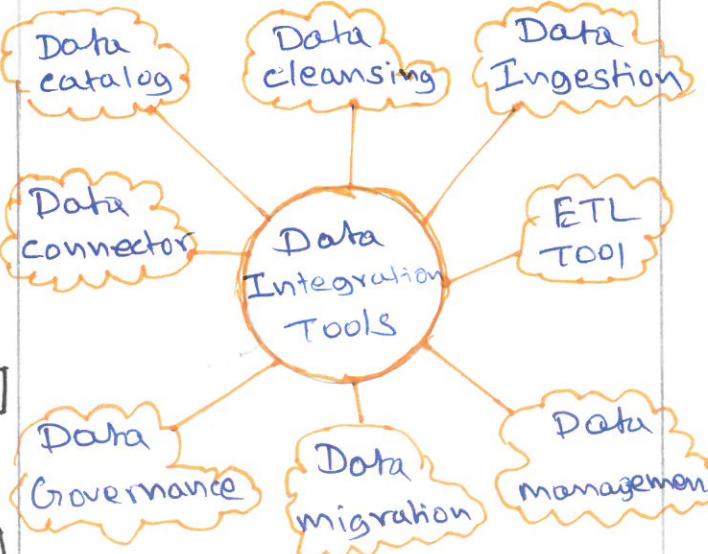
Advantages of Spark



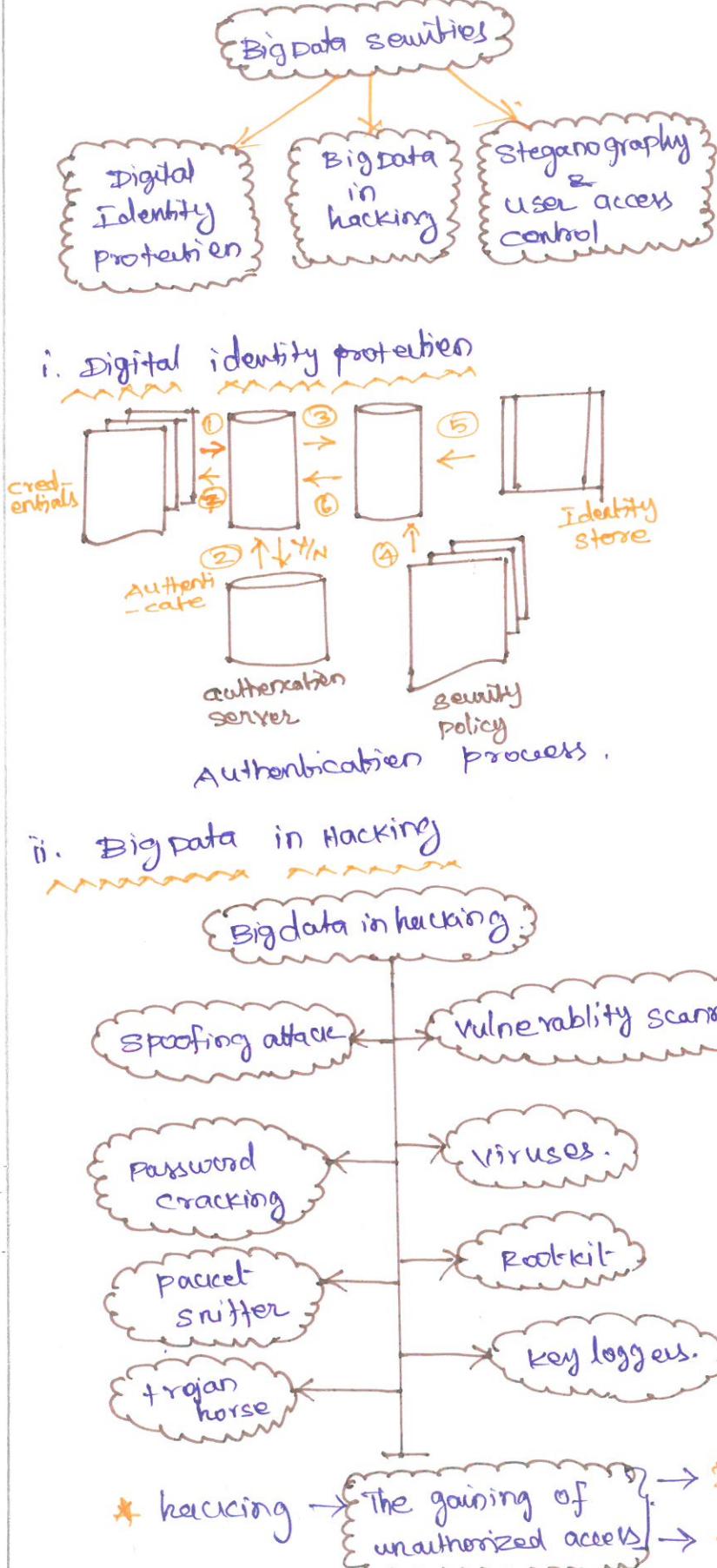
RDD in Spark



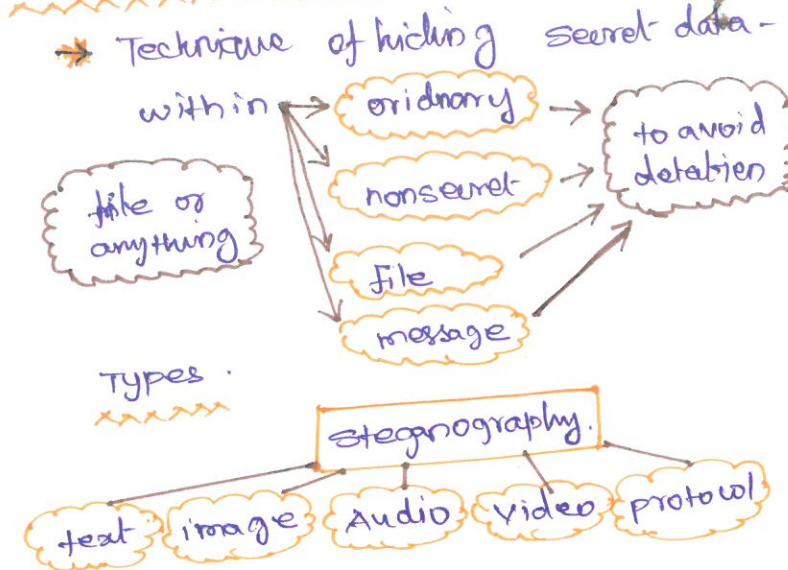
Data Integration Tools



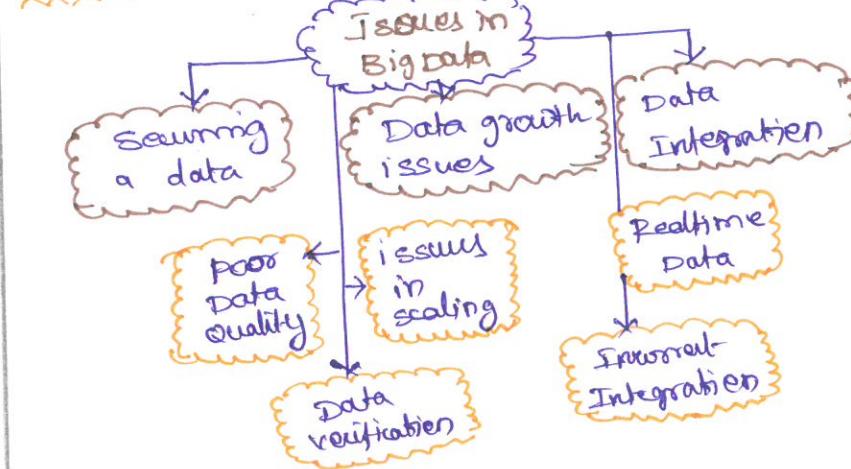
BIG DATA SECURITIES :-



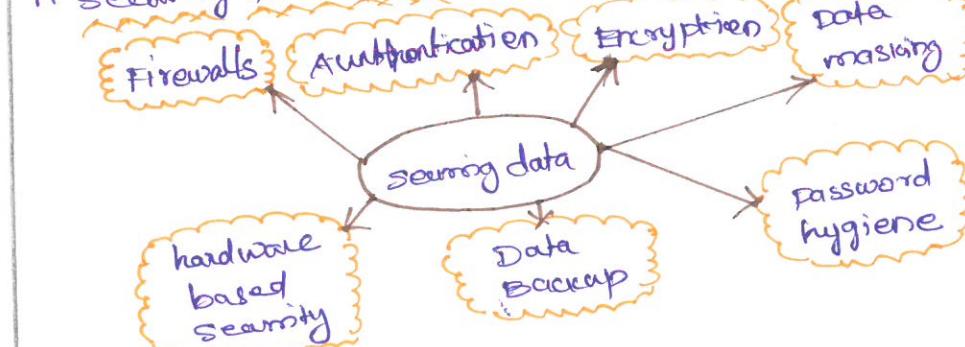
iii. Steganography.



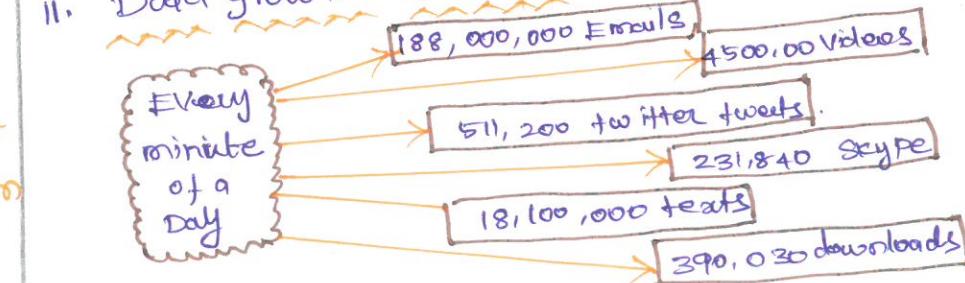
ISSUES IN BIG DATA :-



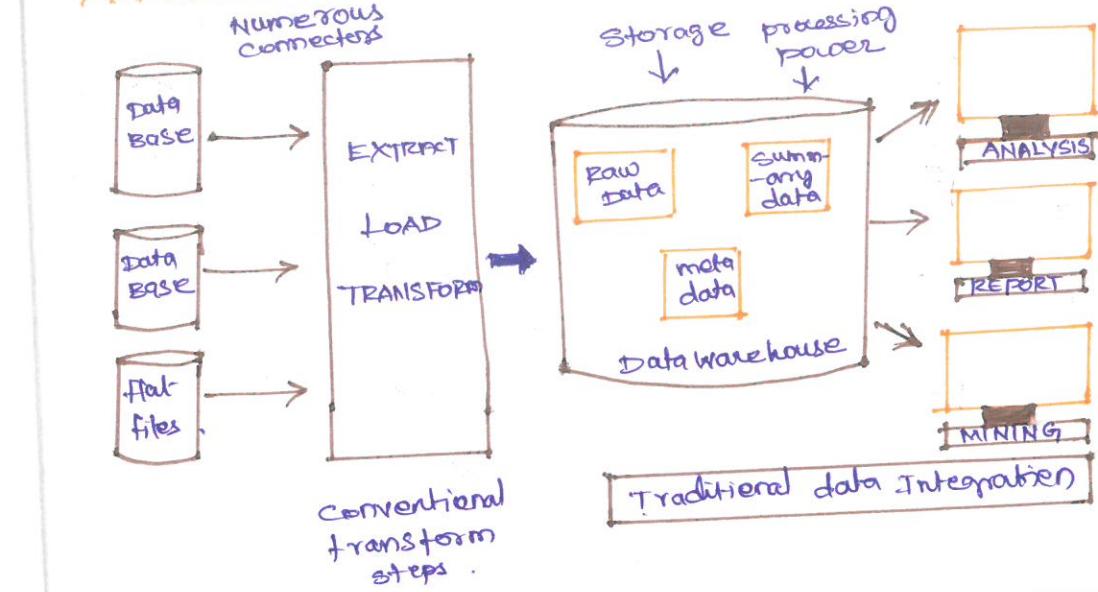
i. Scouring the data :-



ii. Data growth Issues :-

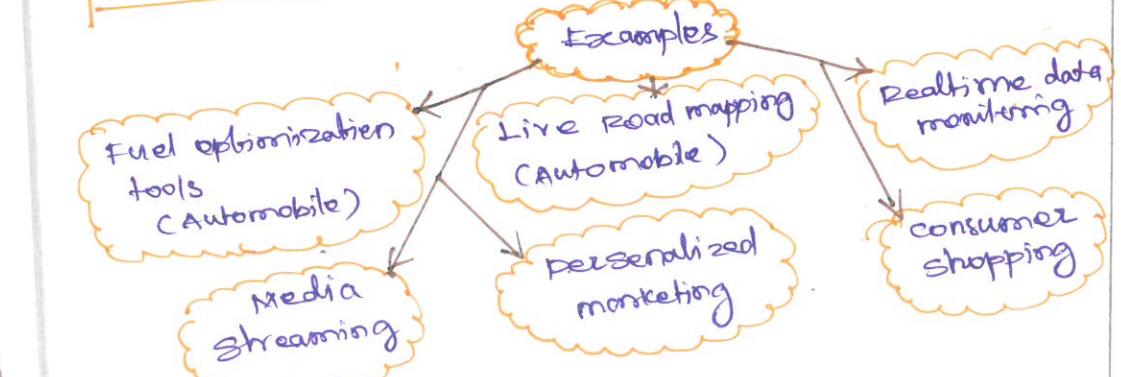
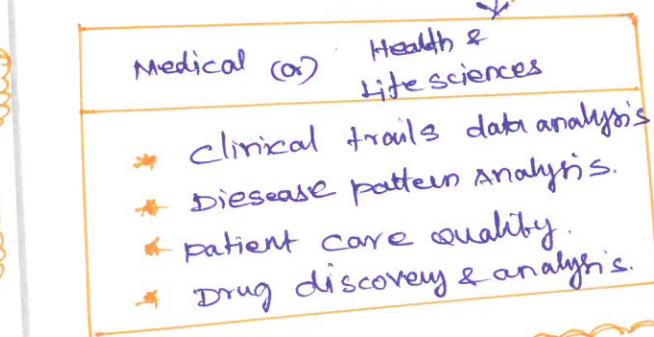


iii. Data Integration :-



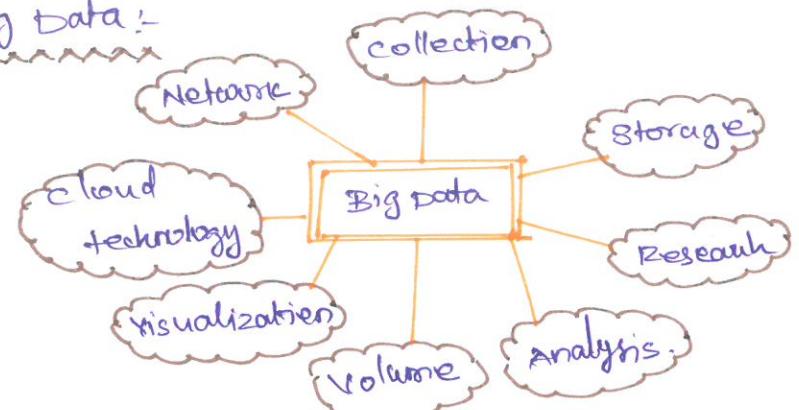
BIG DATA APPLICATIONS :-

- * Applications.
- * Examples.



HADOOP AND MAPREDUCE ARCHITECTURE

Big data:-



Apache Hadoop & Hadoop Ecosystem:-

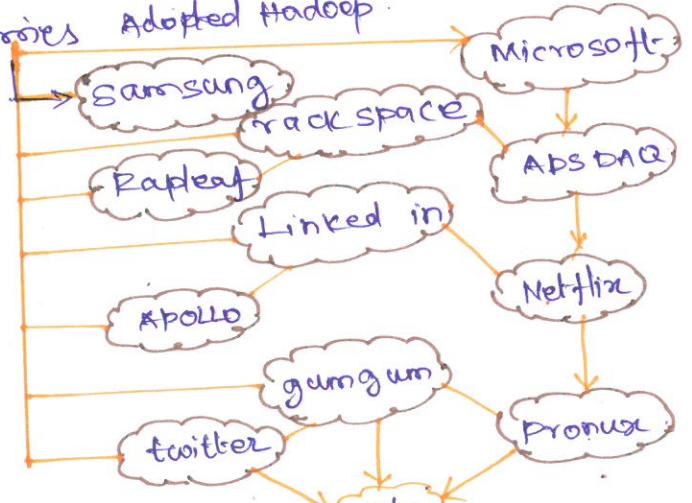
Hadoop:-

scalable fault-tolerant distributed system for
→ Data storage
→ processing

* It operates

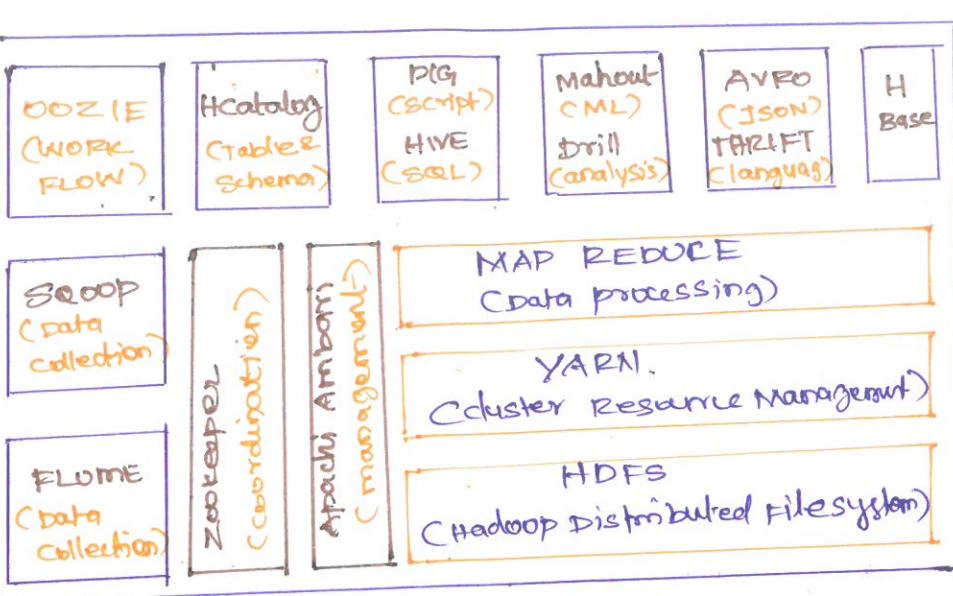
structured data
unstructured data

* Industries Adopted Hadoop

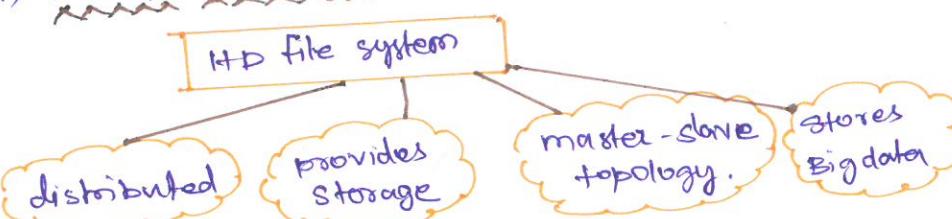


i. Hadoop ECO System:-

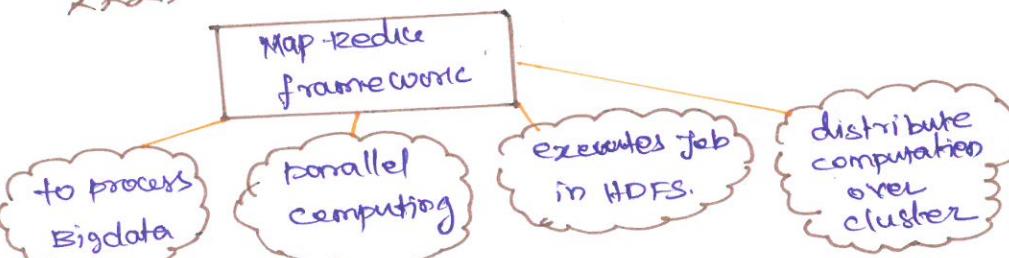
- * Hadoop Distributed file system (HDFS)
- * Map - Reduce (processing layer)
- * YARN - (resource management layer)
- * others (spark, pic, zookeeper ..etc)



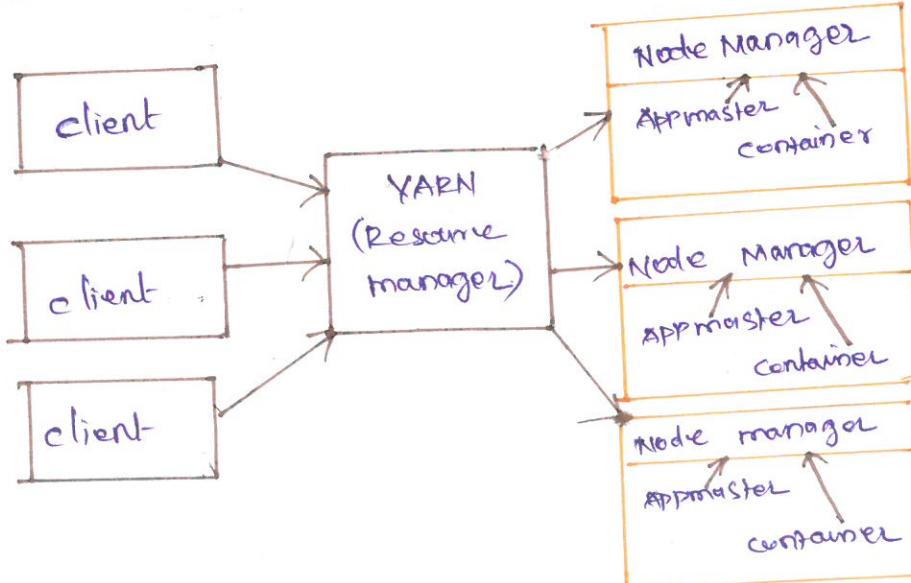
ii) Hadoop distributed File System (HDFS)



iii) Map - Reduce :-



iv) YARN (Yet Another Resource Negotiator)



ANALYZING DATA WITH HADOOP STREAMING.

Hadoop streaming:-

Hadoop streaming is a utility

Enables to create or run Map reduce

script runs in any language

Java / Non-Java as mapper / reducer

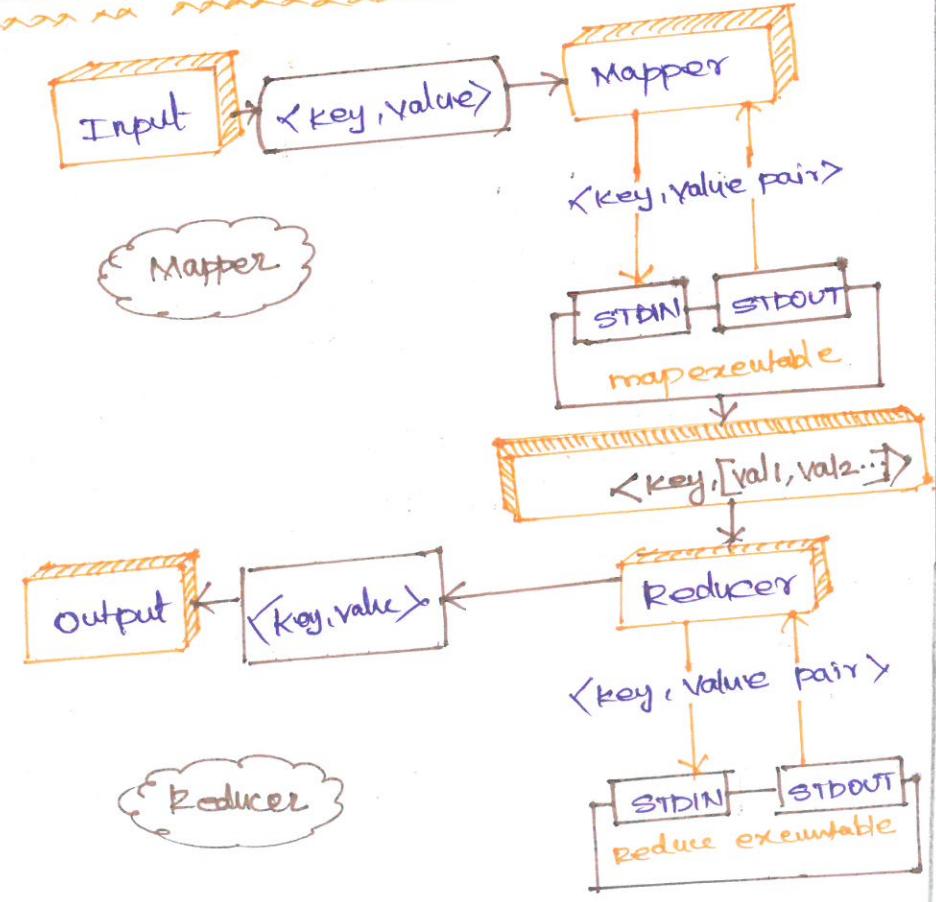
It uses unix streams as an interface.

It reads standard input & writes standard output

Advantages:-

- * Availability
- * Reduced development time
- * faster conversion.
- * performance
- * Learning.

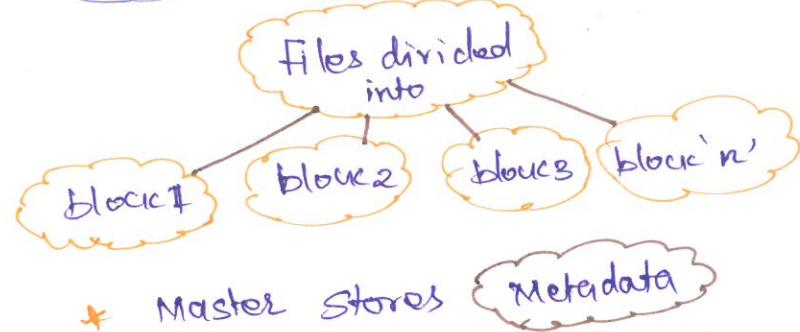
Flow of data Analysis:-



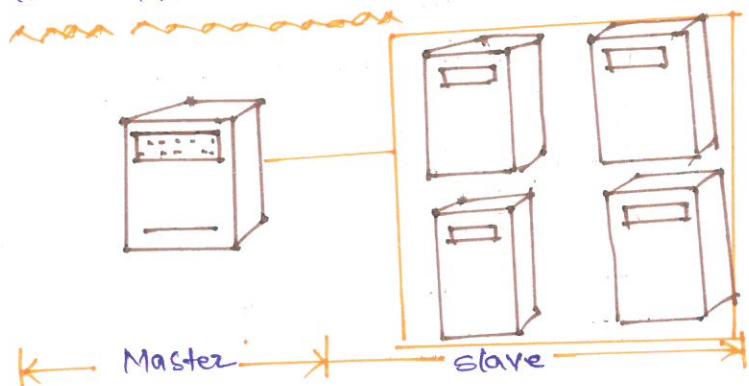
HDFS concepts (Hadoop distributed filesystem)

Introduction:-

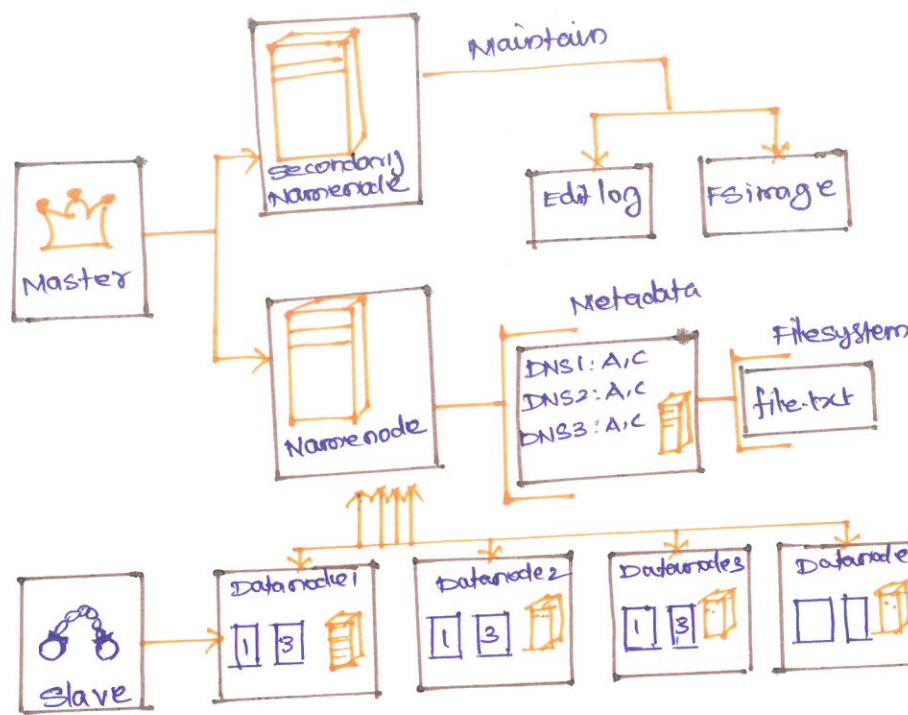
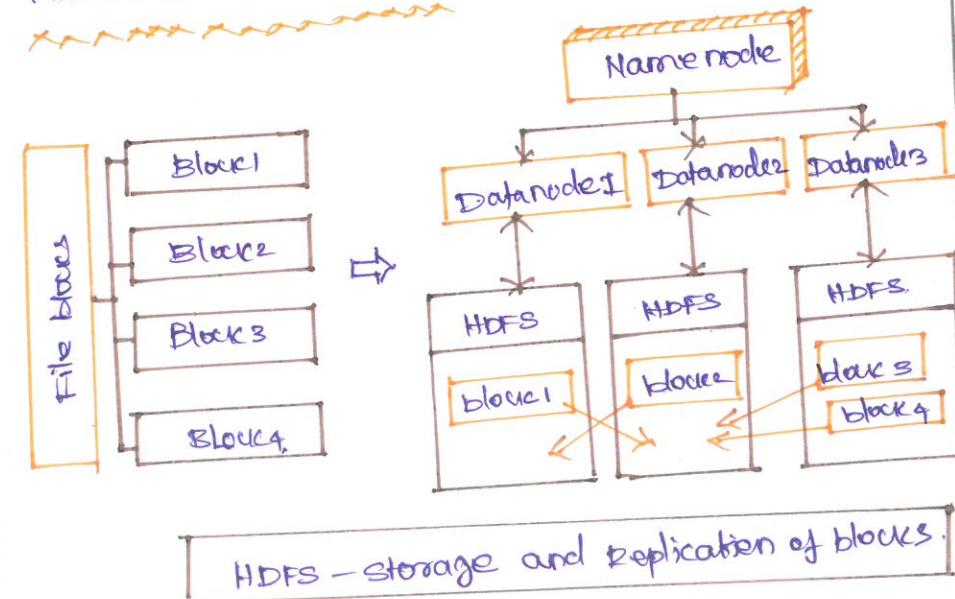
- * provides storage for Hadoop.
- * Master-slave topology.
- Master** → High end machine
- Slave** → Inexpensive computers.



HDFS Architecture:-



WORKING OF HDFS:-

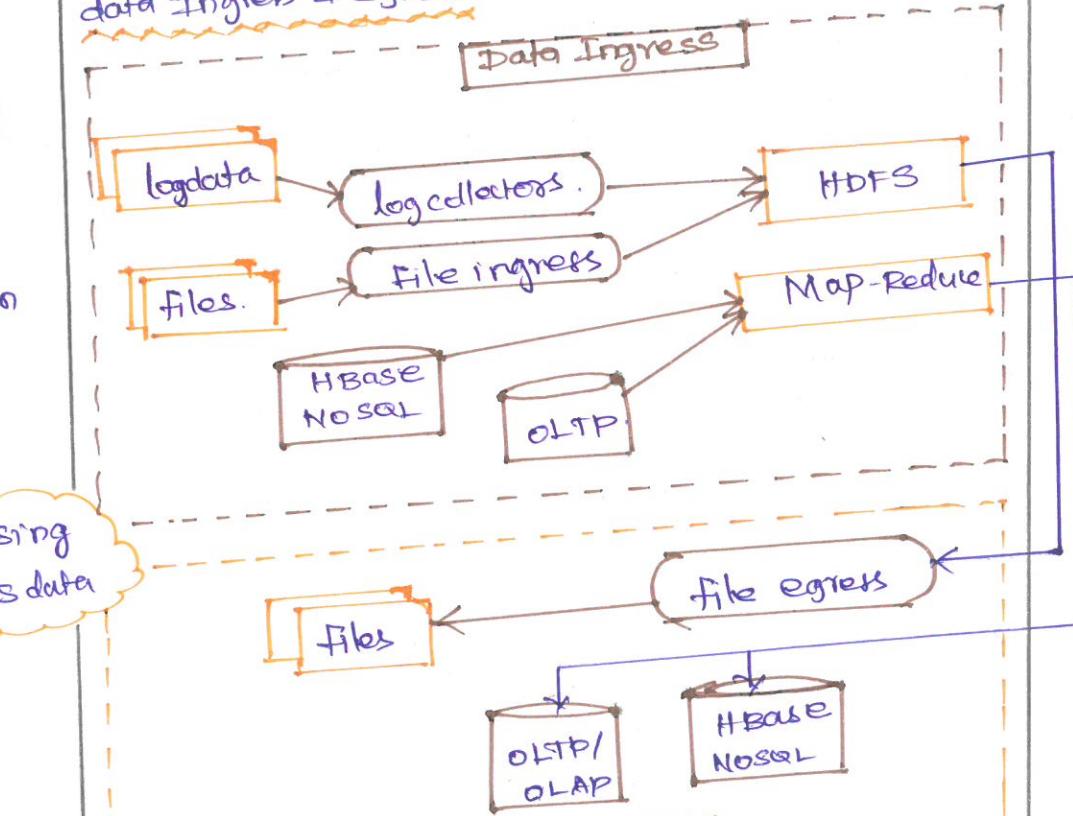
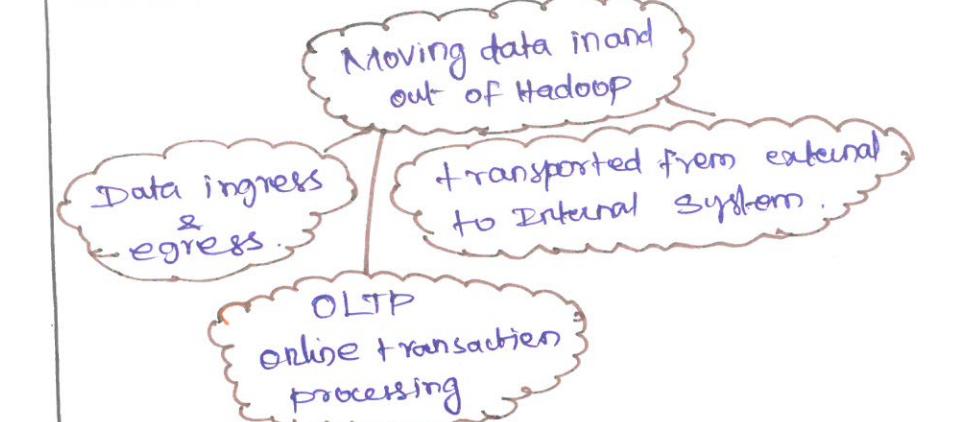


Interface to HDFS:-

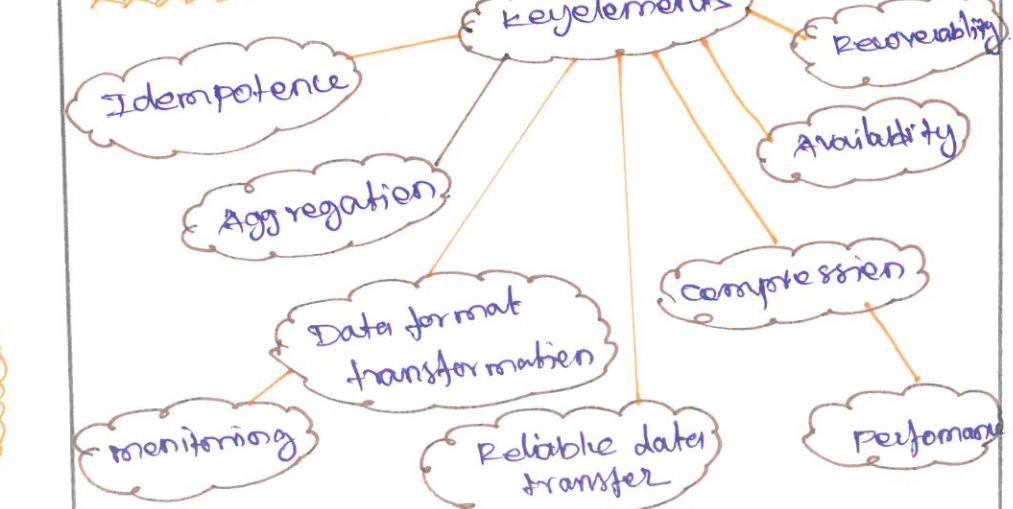
- * Several ways to interact with data stored in HDFS.



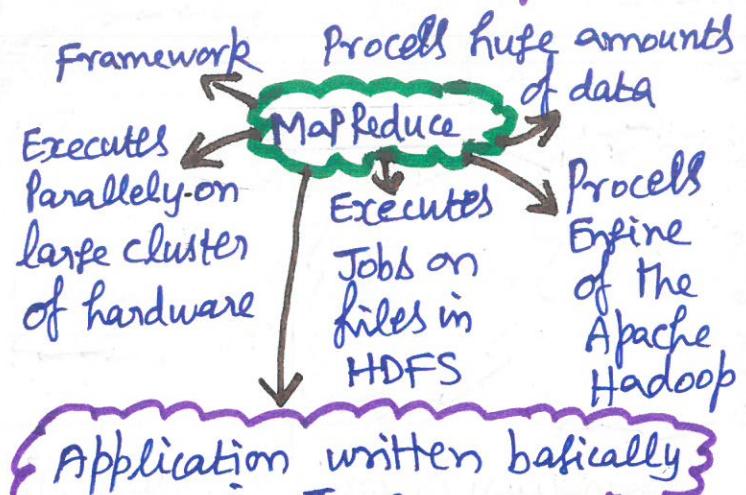
Moving Data in and out of Hadoop:-



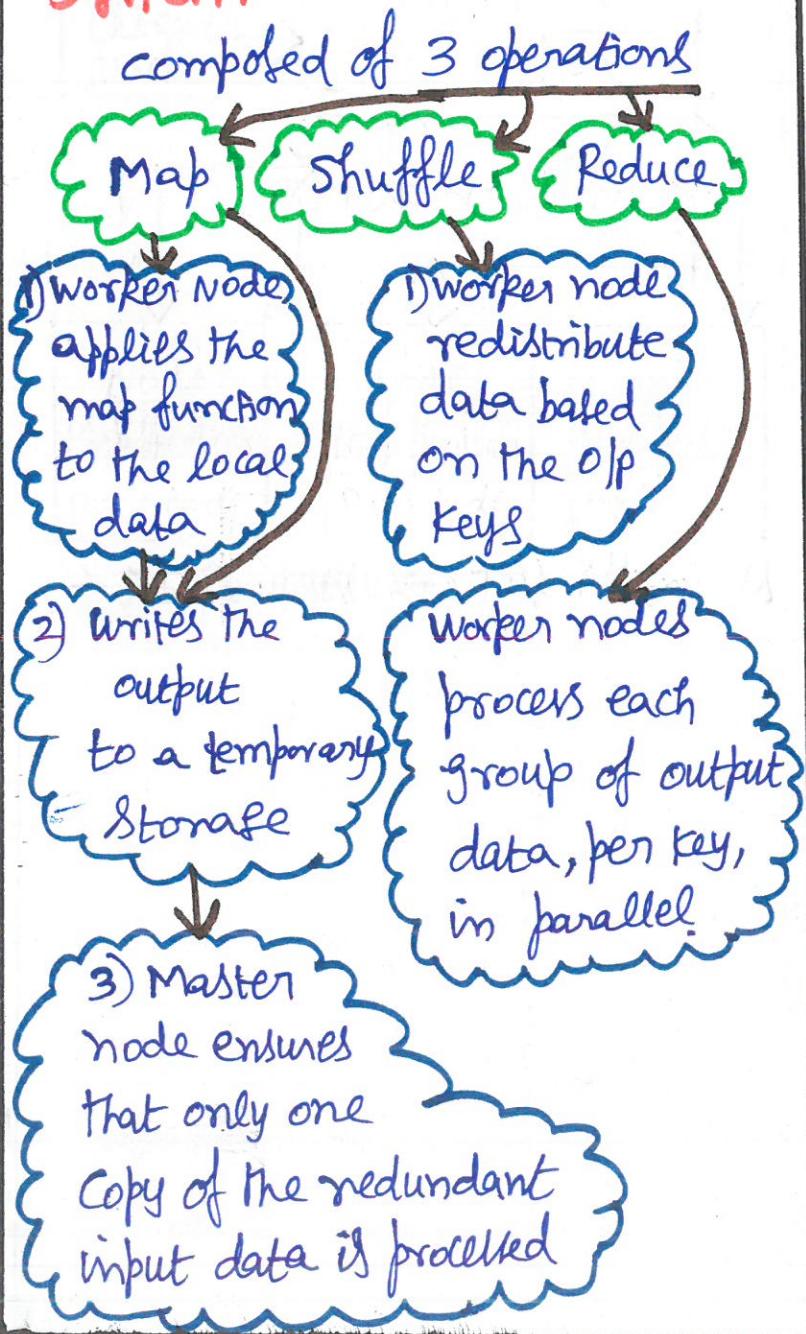
key elements of Data movement:-



Introduction to MapReduce



MapReduce Framework (or) System



Word count Example

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

Reducing

Final Result

Input

Splitting

Mapping

Shuffling

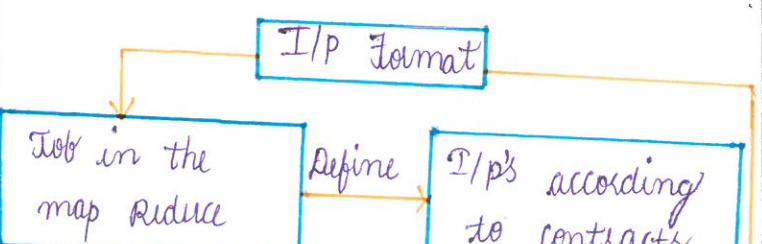
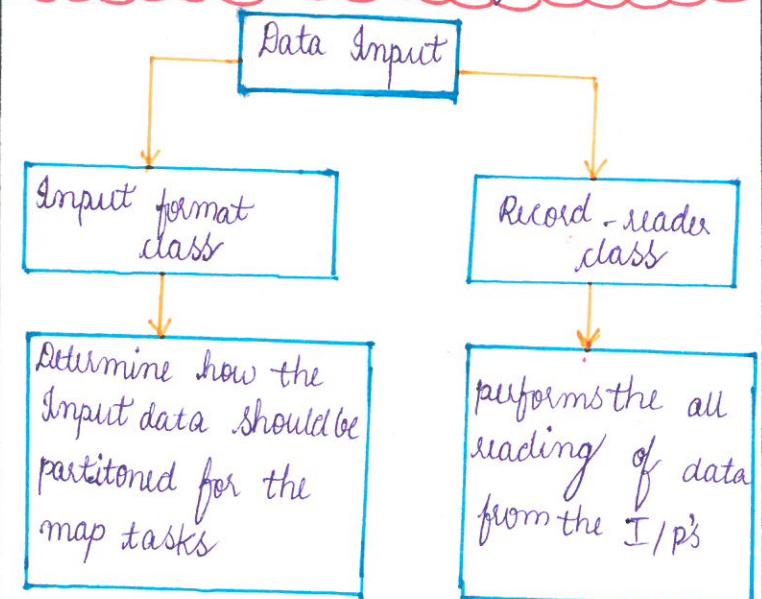
Reducing

</div

Understanding Inputs and outputs of map Reduce:-

Inputs and outputs:-

- * Your Data might be XML files
- * XML (Extensible Markup Language)
- * sits behind a no. of FTP servers.
- * Text log lines of files → central web server
- * Lucene Indexes → HDFS
- * Hadoop Distributed File System

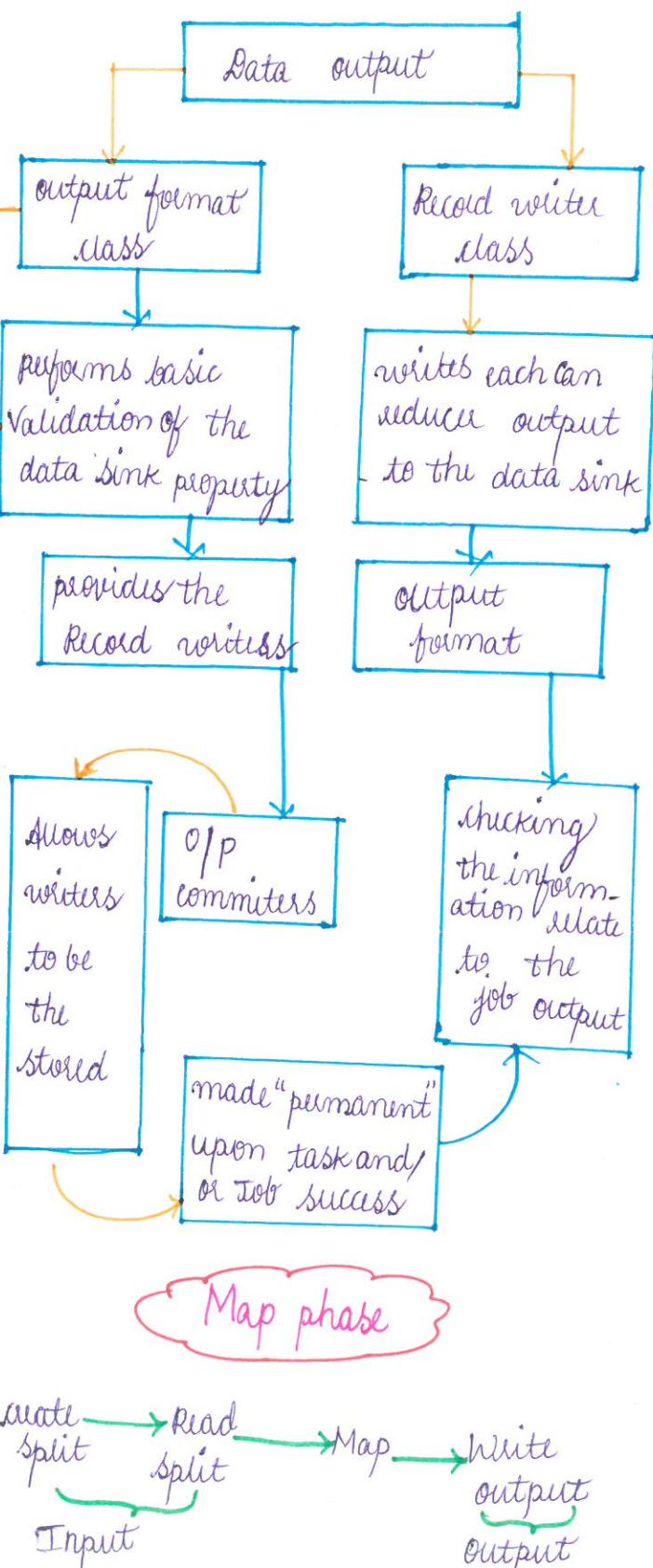


Implementers must fulfill three constraints / contracts

Contact 1: describe type information and for map I/P keys values

Contact 2: specify how the I/p data should be partitioned

Contact 3: Indicate the record reader instance that should read the data from source



Anatomy of Map Reduce Job Run:-

Mapper ,Shuffler & Reducer

Mapper

- 1) perform Map side operation (by you)
- 2) write to in-memory buffer (by framework)

Shuffler

- 1) partitions output key figure out which map goes to which Reduce (by framework)
- 2) several unique key in one partition.
- 3) specify partitioner class which implement partitioner ZD extraction process
- 4) sort first by partitioner ZD then by key value within partition (by framework)
- 5) call every single combine for each key of partition. (by framework)
- 6) spill to disk (also group by key and merge) if limit exceed (Default 100MB).

Reducer

- 1) start read map O/P's from disk
- 2) Merge outputs -

* sort by partition ID and then by key — ③

* Group by key — ④

* call Reduce operation defined by you for each unique key — ⑤

How Map Reduce Job is Running:-

- 1) Map Reduce app submit Job to Hadoop client
- 2) client ask Resource manager to get app ID
- 3) copy job resources to HDFS:
 - a) checks the o/p specification of the job
 - b) computes the I/P splits for the jobs
 - c) copy Jar and throw error if needed.
- 4) submit application
- 5) Resource Manager:
 - a) Allocate container on same node.
 - b) Run application master on that node.
 - c) Application master initialize Job.
 - d) Retrieve Input splits.
 - e) Allocate resources and start of container
 - f) step 10 retrieve resource for map (or) reduce task.
 - g) step 11 run map (or) reduce task.

Failures in classical mapreduce and YARN (Yet Another Resource Negotiator)

Types of Failures

Task Failure

Case 1: Child Task Failing

Map Reduce Task throws a Runtime exception

Solution: The task tracker marks the task attempt as Failed, freeing up a slot to run another task.

Case 2: Sudden exit of the child JVM

Solution: The task tracker notices that the process has exited & marks attempt as failed.

Task Tracker Failure:

- Failure of a Task Tracker is another failure mode
- Task Tracker fails by crashing, or running very slowly, it will stop sending heartbeats to the Job Tracker.

Job Trackers notices Task Tracker that has stopped sending heartbeats remove it from its pool of TaskTrackers to schedule tasks on.

Task Trackers black listed by the Job Tracker

Jobtracker Failure

- Most serious failure mode
- Hadoop has no mechanism to deal with it. Single point of Failure so, Job Fails.
- Low chance of occurring.
- YARN eliminates single point of failure.

Failures in YARN

Task Failure

Due to Runtime exceptions and sudden exits of the JVM

Application master Failure

Due to hardware Failures or Network Failures.

Node Manager Failure

- It will stop sending heartbeats to the resource manager
- Node manager will be removed from the resource manager's pool of available nodes

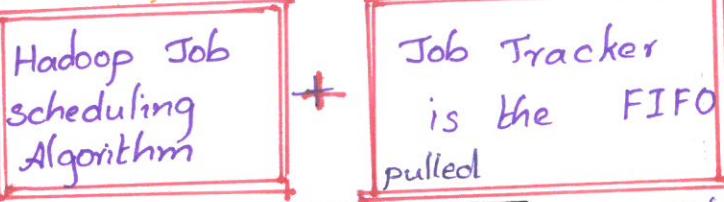
Resource manager Failure

- Neither Jobs nor task container can be launched
- Able to recover from crashes by using checkpointing mechanism to save its state to persistent storage.
- New resource manager instance is brought up and it is recovered.

Job scheduling - Data serialization

Types of Scheduling

1. Hadoop FIFO Scheduler



- No concept of the priority or size of Job

2. Hadoop Fair scheduler



3. Hadoop capacity scheduler

- Multiuser scheduling simulates a separate mapReduce cluster along with FIFO scheduling

Data Serialization

A way of representing data in memory of as a series of bytes

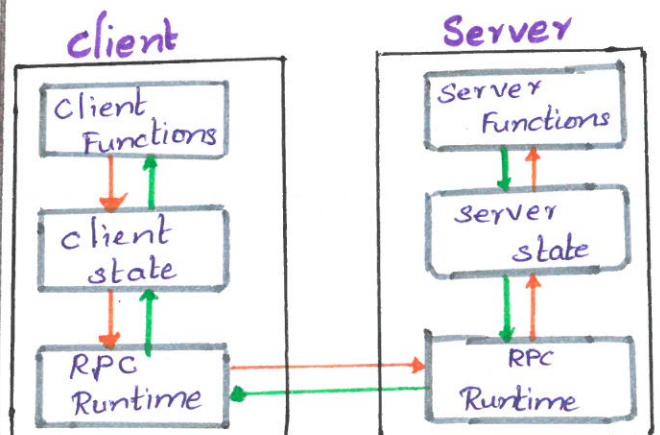
Example : AVRO

- Efficient data serialization framework.
- Widely supported throughout Hadoop and its ecosystem.
- Using Sqoop, data can be imported to HDFS in AVRO and parquet File Format.

Remote procedure call (RPC)

IPC communication

Used for client - server based application
A client has a request message that RPC translates and sends to the server



- Client stub called by the client
- Makes a system call to send messages to the server and puts parameters in the message
- Message sent from C → S by its OS
- Message is passed to the server stub by the server OS
- Parameters removed by server stub
- Server procedure called by server stub.

