

## UNIT-I

According to National Institute of Standards and Technology

### (1) Introduction to cloud computing:

- The term cloud refers to a Server that are accessed over the Internet which is present at remote location.
- cloud can provide services over public and private networks i.e, WAN, LAN (or) VPN
- Applications such as e-mail, web conferencing, customer relationship management (CRM) execute on cloud.
- cloud computing means storing, Managing and accessing the data and programs on remote Servers that are hosted on the Internet instead of computer's hard drive (or) local server.
- cloud computing is also referred to as Internet based computing.
- cloud computing involves provisioning of computing, Networking and storage resources on demand and providing these resources as metered services to the users, in a "Pay as you go" Model.

company - 10 systems we need to see we maintain hardware.

- for that we need to see where we don't need to be maintain Infrastructure

If we take 100-System and 50 employees will left

Resource computers will get waste, All the same time we have some services - will have service.

According to National Institute of Standards and Technology defines cloud computing as:

Definition:-

cloud computing is a model for enabling ubiquitous, convenient, on-demand network access to a shared pool of configurable computing resources (e.g., networks, servers, storage, applications and services) that can be rapidly provisioned and released with minimal management effect (or) service provided interaction

(or) cloud computing is the on-demand availability of computer system resources (especially data storage/cloud storage and computing power) without direct active management by the user.

⇒ Advantages of cloud computing:-

1. Lower cost computer for users
2. Lower IT Infrastructure cost
3. fewer Maintenance cost.
4. Lower software cost.
5. Instant software updates.
6. Increased computing Power
7. unlimited storage capacity
8. cloud computing offers on-demand self-service
9. It offers load balancing that makes it more

o. One can Manipulate and configure the applications  
online at any time.

⇒ Disadvantages of cloud computing:-

1. Require a constant internet connection
2. Stored data might not be secure
3. No control on Resources
4. Depends on cloud providers
5. Interoperability.

(2) characteristics of cloud computing

NIST identifies five essential characteristics of  
cloud computing:

(i) On Demand self Service:- cloud computing allows  
the users to use web services and resource on  
demand, without requiring interactions with the  
cloud service provider (CSP). One can logon to a  
website at anytime & use them.

(ii) Broad Network Access:- The computing services  
are generally provided over standard networks and  
heterogeneous devices/platforms such as workstatio-  
-ns, laptops, tablets and smartphones.

(iii) Resource pooling:- The computing and storage res-  
-ources provided by cloud service providers are

Pool- to serve multiple users.

Multi-tenant aspects of the cloud allow multiple users to be served by the same physical hardware.

(iv) Rapid Elasticity: - cloud computing resources can be provided rapidly and elastically. cloud resources can be rapidly scaled up (or) down based on demand. two types of scaling options exist:

a. Horizontal scaling (Scaling Out): - It involves launching and providing additional server resources.

b. Vertical Scaling (Scaling Up): - It involves changing the computing capacity assigned to the server resources while keeping the number of server resources constant.

(v) Measured Service: - cloud computing resources are provided to users on a "pay-per-use model". The usage of the cloud resources can be monitored, controlled, and reported, providing transparency for both the provider and consumer.

In addition to these five essential characteristics of cloud computing, other common characteristics that again highlight savings in cost include:

(i) Performance: - cloud computing provides improved performance for applications since the resources available to the applications can be scaled up or down based on the dynamic application workload.

(ii) Reduced Costs: - cloud computing provides cost

and storage resources as required can be provisioned dynamically, and upfront investment in purchase of computing assets to cover worst case requirements is avoided. This saves significant cost for organizations and individuals.

Eg:- e-commerce applications

(iii) Outsourced Management: - cloud computing allows the users (individuals, large organizations, small and medium enterprises and governments) to outsource the IT infrastructure requirements to external cloud providers.

(iv) Reliability: - Applications deployed in cloud computing environments generally have a higher reliability since the underlying IT infrastructure is professionally managed by the cloud service. Cloud service providers specify and guarantee the reliability and availability levels for their cloud resources in the form of service level agreements (SLAs).

(v) Multi-tenancy: - The multi-tenanted approach of the cloud allows multiple users to make use of the same shared resources.

- Modern applications such as e-commerce, Business-to-Business, Banking and financial, Retail and social networking applications that are deployed in cloud computing environments are multi-tenanted applications.

- multi-tenancy can be of different forms:

\* Virtual multi-tenancy: - In virtual multi-tenancy, computing and storage resources are shared among multiple users.

\* organic multi-tenancy:- In organic multi-tenancy every component in the system architecture is shared among multiple tenants, including hardware, OS, database servers, application servers, load balancers, etc

### 3. cloud Models:

3.1 - Service Models

3.2 - Deployment Models

#### 3.1 Service Models:

Cloud computing services are offered to users in different forms. NIST defines three cloud service models as follow

a. Infrastructure as a Service (IaaS)

b. Platform as a Service (PaaS)

c. Software as a Service (SaaS)

#### Software as a Service (SaaS)

Applications, management and user interfaces provided over a network

#### Platform as a Service (PaaS)

Application development frameworks, operating systems and deployment frameworks

#### Infrastructure as a Service (IaaS)

Virtual Computing, storage and network resources that can be provisioned on demand.

F. cloud computing service models

a. Infrastructure

→ It provides

→ IaaS is cloud (IaaS)

→ It is a System

→ It is a System

development

→ It provides

such as storage

→ On demand go to

→ we can search

b. Platform

→ PaaS

the platform application

→ The cloud

Operations

→ The deployment

the

c. Software

→ so

→ It

## a. Infrastructure as a Service (IaaS):-

- It provides us infrastructure
- IaaS is also known as Hardware as a Service (HaaS)
- It is a type of cloud computing service used by system administrators/Network architects.
- It simply provides the underlying OS (operating system), security, networking, and servers for developing the applications
- It provides access to fundamental resource such as machines, virtual machines, virtual storage etc.
- On demand, over the internet and on a pay as you go basis
- we can scale up and shrink the resources as per requirement

## b. Platform-as-a Service (PaaS):-

- PaaS cloud computing platform is created for the programmer to develop, test, run, and manage the applications
- The cloud service provider manages the underlying cloud infrastructure including servers, network, operating systems and storage
- The users, themselves, are responsible for developing, deploying, configuring and managing applications on the cloud infrastructure.

## c. Software-as-a-Service (SaaS):-

- SaaS is also known as "on-demand software"
- It is a software in which applications are

→ Users can access the applications with the help of internet connection and web browser

→ maintenance of software and hardware done by the vendor, so it removes the cost of hardware and software maintenance.

### IaaS

<b>Benefits</b> <ul style="list-style-type: none"> <li>• Shift focus from IT management to core activities</li> <li>• pay-per-use / pay-per-go pricing</li> <li>• Dynamic Scaling</li> <li>• Secure access</li> </ul>	<b>characteristics</b> <ul style="list-style-type: none"> <li>- multi tenancy</li> <li>- virtualized hardware</li> <li>- management tools</li> <li>- Disaster recovery</li> </ul>	<b>Example</b> <ul style="list-style-type: none"> <li>- Amazon elastic compute cloud (EC2)</li> <li>- Rackspace</li> <li>- Go Grid</li> <li>- Eucalyptus</li> <li>- Joyent</li> <li>- TerraMark</li> <li>- OpSource</li> </ul>
<b>Adoption</b> <ul style="list-style-type: none"> <li>• Individual users: Low</li> <li>• Small &amp; medium enterprises: medium</li> <li>• Large organizations: High</li> <li>• Government: High</li> </ul>		

### PaaS

<b>Benefits</b> <ul style="list-style-type: none"> <li>• Lower upfront &amp; operations costs</li> <li>• Improved scalability</li> <li>• Higher performance</li> <li>• Secure access</li> </ul>	<b>characteristics</b> <ul style="list-style-type: none"> <li>• multi-tenancy</li> <li>• open integration protocols</li> <li>• App development tools &amp; SDKs</li> <li>• Analytics</li> </ul>	<b>Example</b> <ul style="list-style-type: none"> <li>- Google App Engine</li> <li>- windows Azure Platform</li> <li>- force.com</li> <li>- Rightscale</li> <li>- Heroku</li> <li>- Github</li> <li>- Gigaspace</li> </ul>
<b>Adoption</b> <ul style="list-style-type: none"> <li>• Individual users: Low</li> <li>• Small &amp; medium enterprises: medium</li> <li>• Large organizations: High</li> <li>• Government: Medium</li> </ul>		

### SaaS

<b>Benefits</b> <ul style="list-style-type: none"> <li>• Lower costs</li> <li>• No infrastructure required</li> <li>• Seamless upgrades</li> <li>• Secure</li> <li>• High adoption</li> <li>• On-the-move access</li> </ul>	<b>characteristics</b> <ul style="list-style-type: none"> <li>• multi-tenancy</li> <li>• on-demand software</li> <li>• open integration protocols</li> <li>• social network integration</li> </ul>	<b>Example</b> <ul style="list-style-type: none"> <li>- Google Apps</li> <li>- sales force.com</li> <li>- facebook</li> <li>- zoho</li> <li>- Dropbox</li> <li>- Taleo</li> <li>- microsoft office 365</li> </ul>
<b>Adoption</b> <ul style="list-style-type: none"> <li>• Individual users: High</li> <li>• Small &amp; medium enterprises: high</li> <li>• Large organizations: High</li> <li>• Government: Medium</li> </ul>		

Benefits, characteristics and

3.2 Deployment Models:- NIST also defines four cloud deployment models as follows.

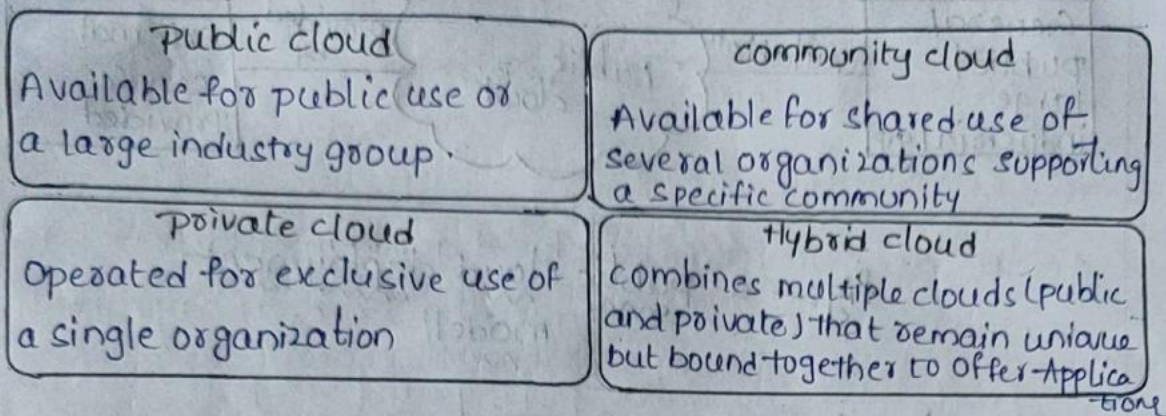


fig:- cloud deployment models

a. public cloud:-

- \* In the public cloud deployment model, cloud services are available to the general public or a large group of companies.
- \* managed by third-party cloud service provider.
- \* public clouds are best suited for users who want to use cloud infrastructure for development and testing of applications and host applications in the cloud to serve large workloads, without upfront investments in IT infrastructure.

Advantages:-

- Minimal Investment
- Infrastructure Management is not required.
- No Maintenance
- Dynamic scalability
- No setup cost

Disadvantages:-

- Less secure because resources are shared

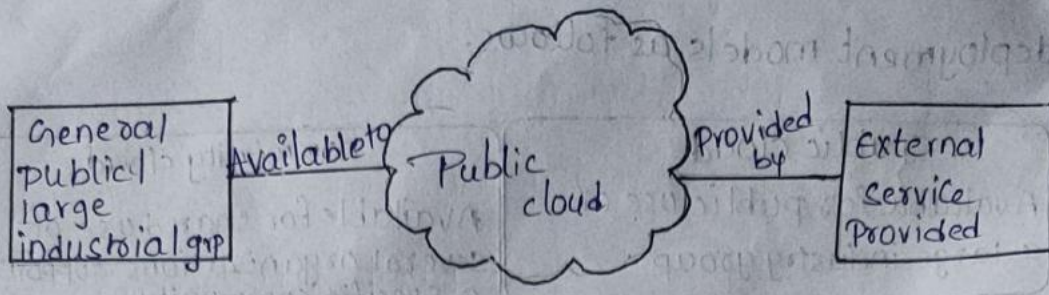


Fig: - a. public cloud deployment model

b. private cloud:-

- \* private cloud services are dedicated for a single organization.
- \* It is managed by Third party (or) the organization internally.
- \* private clouds are best suited for applications where security is very important and organizations that want to have very tight control over their data.

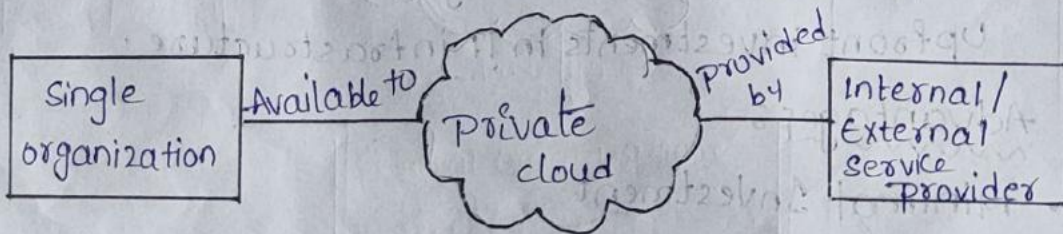


Fig: - b. private cloud Deployment model

Advantages:-

- High security
- Data privacy → only authorized people can access the data.
- More customizable → as companies get to customize their solution as per requirements.

## Disadvantages:-

- private cloud is accessible within an organization so the area of operations is limited.
- High cost → As we need to invest in hardware and software
- Limited scalability

## c. Hybrid cloud:-

- \* The hybrid cloud deployment model combines the services of multiple clouds (private or public)
- \* critical activities are performed by private cloud & non-critical activities by public cloud.
- \* Hybrid clouds are best suited for organizations that want to take advantage of secured application and data hosting on a private cloud, and at the same time benefit from cost savings by hosting shared applications and data in public clouds.

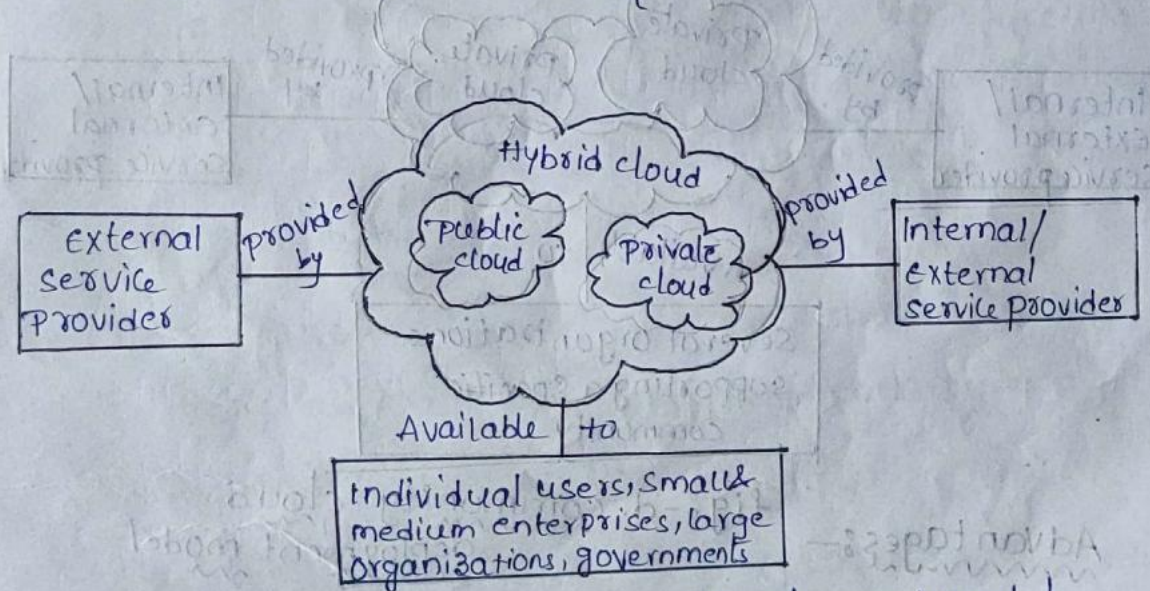


Fig:- c. Hybrid cloud Deployment Model

## Advantages:-

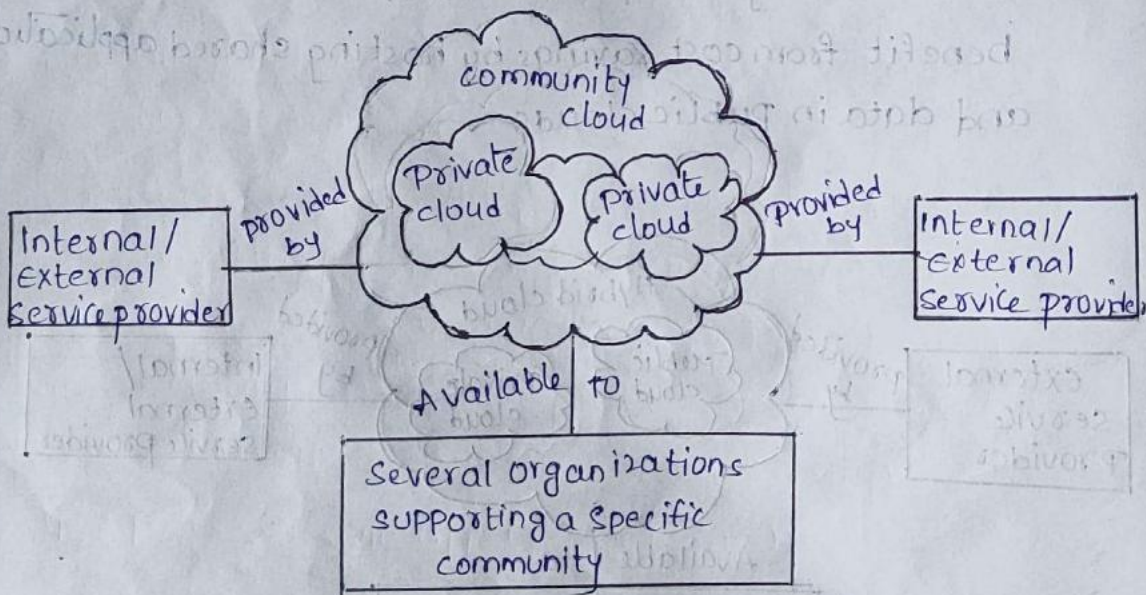
- Scalability
- Security
- Low cost (as compared to private cloud)

## Disadvantages:-

- Managing is complex because there are more than 1 type of deployment model.
- Dependency on infrastructure.

## d. community cloud:-

- \* In the community cloud deployment model, the cloud services are shared by several organizations that have the same policy and compliance considerations.
- \* Available for shared use of several organizations to share the information between the organization and a specific community.



## Advantages:-

- cost effective
- It provides Better Security
- sharing Resources Among companies
- more secured than public cloud but less than private cloud.

## Disadvantages:-

- Data is accessible between organization
- consistent Maintenance cost
- overall cost will be increased compared to public cloud.

## 4. cloud Services Examples

4.1 IaaS: Amazon EC2, Google compute Engine, Azure Vms

4.2 PaaS: Google App Engine

4.3 SaaS: Salesforce

4.1 IaaS: Amazon EC2, Google compute Engine, Azure Vms

(1) AWS is an Infrastructure-as-a-Service (IaaS) offered from Amazon.com

(2) Elastic compute cloud (EC2) is a Webservice that provides computing capacity in the form of virtual machines that are launched in Amazon's cloud computing environment.

(3) It is used to launch as many (or) a few virtual servers as you need, configure security and networking, and manage storage.

(4) It is enables you to scale up (or) down to handle changes in requirements, reducing your need to forecast traffic.

(5) Amazon provides pre-configured Amazon Machine Images (AMIs) which are templates of cloud

(6) The pricing model for EC2 instance is based on a pay-per-use model

(7) AEC2 provides number of powerful features for building scalable and reliable applications such as auto-scaling and elastic load balancing.

(8) users can also create their own AMIs with custom Applications, libraries and data.

(9) Google Compute Engine (GCE) is a IaaS offered from Google.

(10) It provides Virtual Machines of various computing capacities ranging from small instances (Eg: 1 virtual core with 1.38 GCE unit and 1.7 GB memory) to high memory machine types (Eg: 8 virtual core with 32 GCE units and 52 GB memory).

(11) windows Azure Virtual machines is an IaaS offered from Microsoft.

(12) Azure Vms provides virtual machines of various computing capacities ranging from small instances (1 virtual core (vc) with 1.75 GB memory) to memory intensive machine types (8 virtual cores with 56 GB memory).

4.2 paas: Google App Engine (GAE)

(1) GAE is a platform-as-a service offered from Google

(2) It is a cloud-based service for hosting web applications and storing data.

(3) It allows users to build scalable and reliable applications that run on the same systems that power Google's own Applications.

- (5) Developers can develop and test their applications with GAE SDK on a local machine and then upload it to Google App Engine (GAE) with a simple click of a button.
- (6) Applications hosted in GAE are easy to build, maintain and scale.
- (7) GAE supports applications written in several programming languages.
- (8) With GAE's Java runtime environment developers can build applications written in several programming languages.
- (9) GAE also provides runtime environment for Python languages.
- (10) The pricing model for GAE is based on the amount of computing resources used.
- (11) GAE provides free computing resources for applications upto a certain limit. Beyond that limit, users are billed based on amount of resources used.

#### 4.3. SaaS: Salesforce

- (1) Salesforce is a cloud-based customer relationship management software offered from software-as-a-service.
- (2) Users can access CRM application from anywhere through internet-enabled devices such as workstations, laptops, tablets and smartphones.
- (3) Salesforce is a popular CRM tool for support, sales and marketing teams worldwide.
- (4) Salesforce services allow businesses to use cloud technology to better connect with customers.

(15) For Example: LinkedIn - It brings companies and customers together on the CRM.

(16) sales force allows customer representatives to manage customer profiles, track opportunities, optimize campaigns from lead to close and monitor the impact of campaigns.

(17) Sales force Service cloud is a cloud-based customer service management SaaS.

(18) Service cloud provides companies a call-center like view and allows creating, tracking, routing and escalating cases, also provides self service capabilities to customers.

(19) Salesforce marketing cloud is a cloud based social marketing SaaS.

(20) Marketing cloud allows companies to identify sales from social media, identify the most trending information on any topic and allows companies to actively engage with customers, to create and deploy social content, Manage and execute optimized social advertisement campaigns and track the performance of social campaigns.

(21) Some of the tools included in the sales force sales, service and marketing clouds include:

- Accounts and contacts
- Chatter
- Leads
- Analytics and

# Cloud Based Services and applications

## 5.1 Cloud Computing for Healthcare:

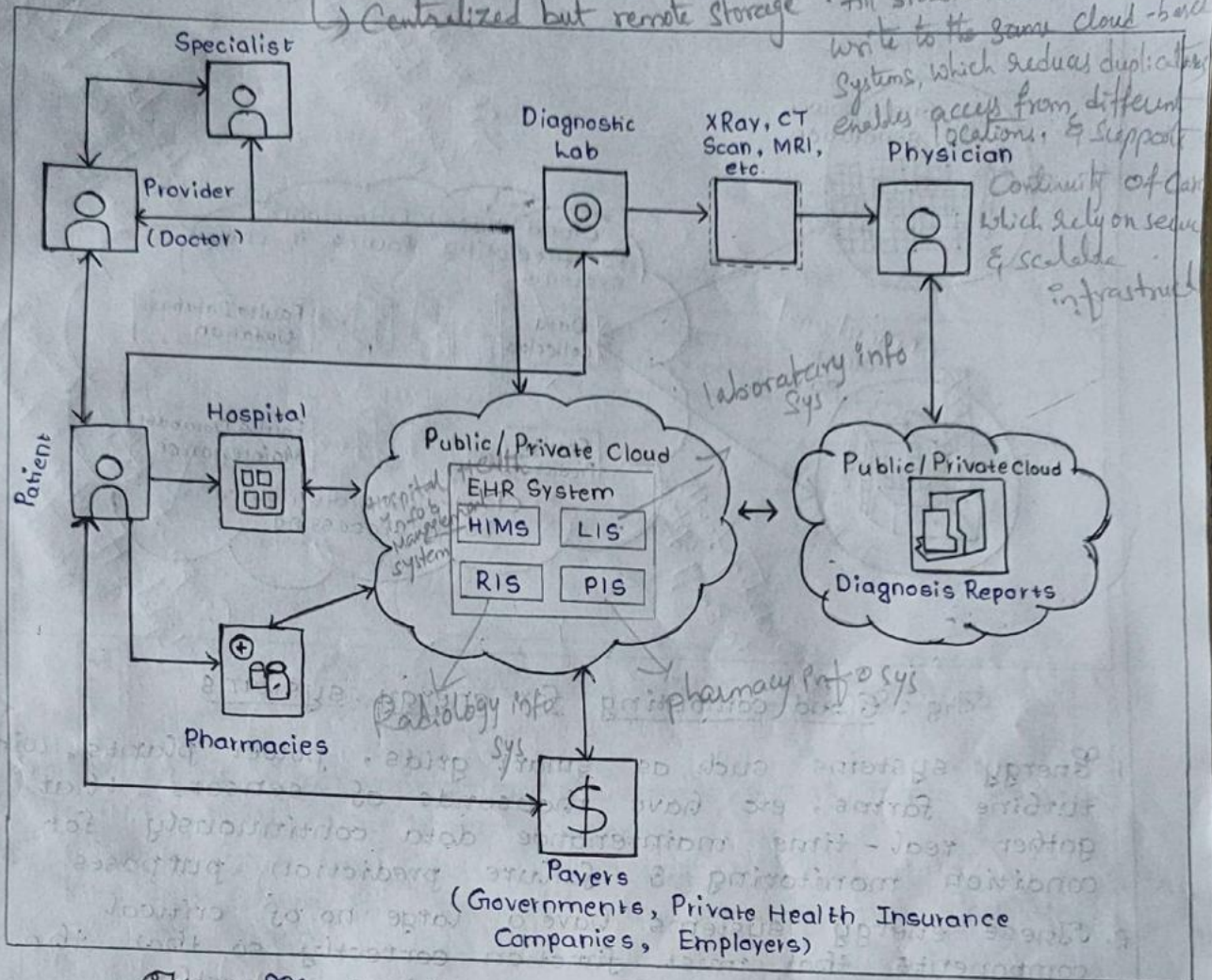


Fig: Cloud computing for Healthcare

1. The above figure shows the application of cloud computing environments to the healthcare system.
2. Hospitals and their affiliated providers can securely access patient data stored in the cloud and share the data with other hospitals and physicians.
3. Patients can access their own health information from all of their care providers and store it in a personal health record (PHR).
4. The PHR can be a vehicle for e-prescribing, a technique known to reduce paper prescriptions and allows physicians and other medical practitioners to write & send prescriptions to a participating pharmacy electronically.
5. History and information staged stored in the cloud using SaaS applications can streamline the admissions, care & discharge processes by eliminating redundant data collection and entry.
6. Health payers can increase the effectiveness & lower the

## Distributed Computing With Hadoop And Spark In Cloud Computing For AI.

Apache Hadoop is an open-source software framework for storing and processing massive datasets across clusters of inexpensive computers. It is designed for large-scale, fault-tolerant, **batch processing** of data stored on its distributed file system (HDFS). Apache Spark is an open-source, general-purpose cluster-computing framework built for speed and sophisticated analytics. It is designed to overcome the limitations of Hadoop's MapReduce by using **in-memory processing**, making it up to 100 times faster for certain workloads, particularly for real-time analytics and machine learning.

In 2025, distributed computing using **Apache Hadoop** and **Apache Spark** remains the foundational architecture for large-scale Artificial Intelligence (AI) and machine learning (ML) in the cloud. While Hadoop provides a cost-effective, scalable storage backbone, Spark acts as the high-speed engine designed specifically for the iterative nature of AI workloads.

### Core Roles in AI Workflows

Organizations often use these two frameworks together in a hybrid model to maximize efficiency:

- **Hadoop (The Foundation):** Best for storing petabytes of raw, unstructured data in its **Hadoop Distributed File System (HDFS)**. In AI, it is primarily used for **data lakes**, long-term archiving, and heavy batch-oriented ETL (Extract, Transform, Load) tasks that prepare data for model training.
- **Spark (The Engine):** Specialized for high-speed processing through **in-memory computing**, which can be up to 100x faster than Hadoop's MapReduce for iterative AI algorithms. Its **MLlib** library provides built-in, scalable machine learning algorithms for classification, regression, and clustering.

### Distributed Computing for AI Tasks

AI Requirement	Framework/Component	Why It's Used
Model Training	Spark MLlib	Handles the iterative math required for ML by keeping data in RAM instead of reading from disk constantly.
Real-time Inference	Spark Streaming	Enables processing of live data streams (e.g., fraud detection or sensor data) for

1/10  
21

2

		immediate AI insights.
<b>Data Preparation</b>	<b>Hadoop MapReduce</b>	Efficiently cleans and transforms massive historical datasets that are too large for memory.
<b>Graph Analytics</b>	<b>Spark GraphX</b>	Models complex relationships, such as social network connections or recommendation links, using parallel graph processing.

### Cloud Integration in 2025

Cloud platforms provide managed versions of these frameworks, removing the need for manual infrastructure management:

- **AWS (Amazon EMR):** Automatically scales Spark and Hadoop clusters and integrates with S3 for storage.
- **Google Cloud (Dataproc):** Offers fast provisioning (often under 90 seconds) and deep integration with **Vertex AI** for end-to-end ML workflows.
- **Microsoft Azure (HDInsight):** Provides managed clusters with enterprise-grade security and integration with Azure Data Lake.
- **Databricks:** A unified, cloud-native platform built on Spark that is optimized for collaborative data science and AI.

### Key Advantages for AI

- **Fault Tolerance:** Both systems are designed to handle hardware failures automatically. Hadoop replicates data across nodes, while Spark uses **Resilient Distributed Datasets (RDDs)** to recompute lost data if a node fails.
- **Scalability:** Both frameworks scale horizontally, meaning you can process more data or train more complex models simply by adding more commodity servers in the cloud.
- **Polyglot Support:** Spark supports high-level APIs in **Python (PySpark)**, Scala, R, and Java, making it accessible to data scientists who prefer Python for ML development.

OR

In cloud computing environments for AI, **Hadoop provides scalable and cost-effective distributed data storage**, while **Spark offers a fast, in-memory processing engine with built-in machine learning libraries (MLlib)** for distributed

AI workloads. The two frameworks are often used in a hybrid architecture to leverage their respective strengths.

### **Hadoop in the Cloud for AI**

Hadoop primarily serves as a foundational layer for managing massive, diverse datasets, which are essential for training AI models.

- **Distributed Storage (HDFS):** The Hadoop Distributed File System (HDFS) stores vast amounts of structured and unstructured data across a cluster of commodity machines. In the cloud, this data is often stored on scalable object storage services like Amazon S3, Google Cloud Storage, or Azure Blob Storage, with Hadoop providing the distributed file system abstraction.
- **Cost-Effectiveness:** Hadoop is ideal for low-cost, long-term storage and large-scale batch processing, such as initial data ingestion and ETL (Extract, Transform, Load) pipelines, where speed is not the primary concern.
- **Scalability and Fault Tolerance:** It scales horizontally by adding more nodes and provides high fault tolerance through data replication, ensuring data availability for AI model training even if individual nodes fail.

### **Spark in the Cloud for AI**

Spark is a powerful, general-purpose processing engine designed for speed and versatility, making it particularly well-suited for the iterative nature of AI and machine learning (ML) tasks.

- **In-Memory Processing:** Spark processes data in RAM, which makes it significantly faster than Hadoop's disk-based MapReduce for iterative algorithms common in ML model training.
- **Machine Learning Library (MLlib):** Spark includes a built-in, scalable ML library (MLlib) that provides various algorithms for classification, regression, clustering, and more. This allows data scientists to build and deploy ML models efficiently on large datasets.
- **Real-Time Analytics:** Spark Streaming and Structured Streaming enable the processing of live data streams, which is crucial for AI applications like real-time fraud detection or recommendation engines that require immediate insights.

**Ease of Use:** Spark offers high-level APIs in popular languages like Python (PySpark), Scala, R, and Java, making it more accessible to data scientists and developers.

### Hybrid Architecture in the Cloud

In a typical cloud AI architecture, Hadoop and Spark work in concert:

- **Hadoop (Storage Layer):** Stores the raw, massive datasets in a data lake environment.
- **Spark (Compute Layer):** Reads data from the Hadoop storage layer (or cloud object storage), processes it in-memory, and runs iterative ML algorithms using MLlib. The processed data or trained models are then written back to storage.

Cloud platforms such as AWS (Amazon EMR), Google Cloud (Dataproc), and Microsoft Azure offer managed services that simplify the deployment and scaling of both Hadoop and Spark clusters, allowing organizations to focus on AI development rather than infrastructure management.

## 2. What is data ingestion?

Data ingestion is the process of collecting raw data from diverse sources (like apps, IoT devices, databases) and moving it into a central system (data warehouse, lake, database) for storage, processing, and analysis, forming the crucial first step in any data pipeline to enable business intelligence, AI, and decision-making. It involves getting data from many places into one accessible spot, handling different formats and speeds (batch vs. real-time), and making it ready for use by downstream tools.

Data ingestion and processing pipelines in the cloud automate moving data from sources (databases, IoT, APIs) to storage (lakes, warehouses) and then transforming it for analysis, using cloud services for scalability, elasticity, and cost-efficiency, primarily through batch processing (scheduled) or real-time streaming (as it arrives). These pipelines are foundational for big data, enabling business intelligence, machine learning, and quick decision-making by handling diverse data volumes and speeds efficiently.

Data ingestion pipeline is a crucial component of modern data architecture, enabling businesses to efficiently manage and utilize their data. It's the process of importing, transferring, loading, and processing data for later use or storage in a database. This process is integral to data systems, as it's the first step in making raw data accessible and usable for analytics and decision-making.

## What is data ingestion?

Data ingestion is the process of collecting, importing, and processing data from various sources for storage in a database. This process is the first step in making raw data accessible and usable for analytics and decision-making. The data sources

5

can be numerous and diverse, including databases, servers, other data centers, and even online sources.

### **Understanding the data ingestion process**

Data ingestion involves taking data from various sources and moving it to a system where it can be stored and analyzed. The data sources can be numerous and diverse, including databases, servers, other data centers, and even online sources. The data is then processed and loaded into a destination system, such as a data warehouse or a data lake, where it can be accessed and analyzed by data scientists and other users.

The data ingestion process can be done in real-time or in batches. Real-time data ingestion involves importing data as soon as it is produced, providing users with up-to-date, real-time insights. Batch data ingestion, on the other hand, involves collecting and importing data at regular intervals, whether it's every few hours, once a day, or once a week.

### **What are the different types of data ingestion?**

There are primarily two types of data ingestion methods: real-time and batch.

Real-time data ingestion involves importing data as soon as it is produced, providing users with up-to-date, real-time insights. This method is crucial for businesses that need to make immediate decisions based on the latest data.

Batch data ingestion, on the other hand, involves collecting and importing data at regular intervals, whether it's every few hours, once a day, or once a week. This method is suitable for businesses that do not require real-time data and can make decisions based on periodic data updates.

### **The importance of data ingestion pipelines**

Data ingestion pipelines are essential for efficient data management. They automate the data ingestion process, ensuring that data from various sources is consistently and accurately imported, processed, and stored. This not only saves time and resources but also ensures that the data is reliable and ready for analysis.

Moreover, data ingestion pipelines are crucial for businesses that rely on real-time data. They enable these businesses to make timely decisions based on the most recent data. For instance, businesses can use real-time data ingestion pipelines to monitor their website traffic and make immediate adjustments to their marketing strategies based on the incoming data.

### **What are the benefits of data ingestion?**

Data ingestion offers several benefits:

1. **Efficiency:** Automated data ingestion pipelines save time and resources by streamlining the process of importing, processing, and storing data.

- (2)
2. **Reliability:** These pipelines ensure that the data is consistently and accurately imported, making it reliable for analysis.
  3. **Timely Insights:** Real-time data ingestion pipelines provide up-to-date insights, enabling businesses to make timely decisions.
  4. **Scalability:** As businesses grow, so does the amount of data they generate. Data ingestion pipelines can scale to handle increasing data volumes.

### **Building an effective data ingestion pipeline**

Building an effective data ingestion pipeline involves several key steps:

1. **Identify the Data Sources:** The first step in building a data ingestion pipeline is to identify the data sources. These could be databases, online sources, servers, or other data centers.
2. **Determine the Destination System:** The next step is to determine where the data will be stored and analyzed. This could be a data warehouse, a data lake, or another type of data storage system.
3. **Choose the Data Ingestion Method:** The data ingestion method could be either real-time or batch, depending on the needs of the business.
4. **Design the Data Ingestion Process:** This involves designing the process of importing, processing, and loading the data into the destination system. This process should be automated to ensure consistency and accuracy.
5. **Monitor and Optimize the Data Ingestion Pipeline:** Once the data ingestion pipeline is in place, it's important to monitor its performance and make any necessary adjustments to ensure it's running efficiently.

### **Common data ingestion challenges**

Despite its benefits, data ingestion can present several challenges:

1. **Data Variety:** With numerous data sources, managing different data types can be complex.
2. **Data Volume:** As businesses grow, so does the volume of data, making it challenging to manage and process.
3. **Data Velocity:** The speed at which data is generated and processed can be overwhelming, especially for real-time data ingestion.
4. **Data Veracity:** Ensuring the accuracy and reliability of data is crucial, as poor data quality can lead to incorrect insights and decisions.

### **Data ingestion and ELT**

Extract, Load, Transform (ELT) is a modern approach to data integration that focuses on faster data availability, flexibility, and scalability. In an ELT pipeline, data is extracted from sources, loaded into storage, and then transformed as needed. This approach decouples the



extraction and transformation processes, allowing for quicker data ingestion and giving analysts and data scientists faster access to data.

### Data ingestion approaches

There are two main approaches to data ingestion: manual coding and data integration platforms.

Manual coding involves engineers and developers writing each line of code required to build a data pipeline. This approach is time-consuming and labor-intensive.

Data integration platforms, on the other hand, provide pre-built connectors and transformations to streamline the data pipeline process. These platforms are fully managed, reducing the need for constant updates and maintenance.

(or)

### Quick Takeaways

- *Data ingestion is how new data is absorbed into a system.*
- *A data ingestion pipeline is how data is moved from its original sources to centralized storage.*
- *A data ingestion framework determines how data from various sources is ingested into a pipeline.*
- *The data process flow describes how data moves into and through the data pipeline.*

### What is Data Ingestion? An Easy Explanation

**Data ingestion** is the process of collecting and importing data files from various sources into a database for storage, processing, and analysis. The goal of data ingestion is to clean and store data in an accessible and consistent central repository, preparing it for use within the organization.

An enterprise system may use data ingestion tools to import data from dozens or even hundreds of individual sources, including:

- Internal databases
- External databases
- SaaS applications
- CRM systems
- Internet of Things sensors
- Social media

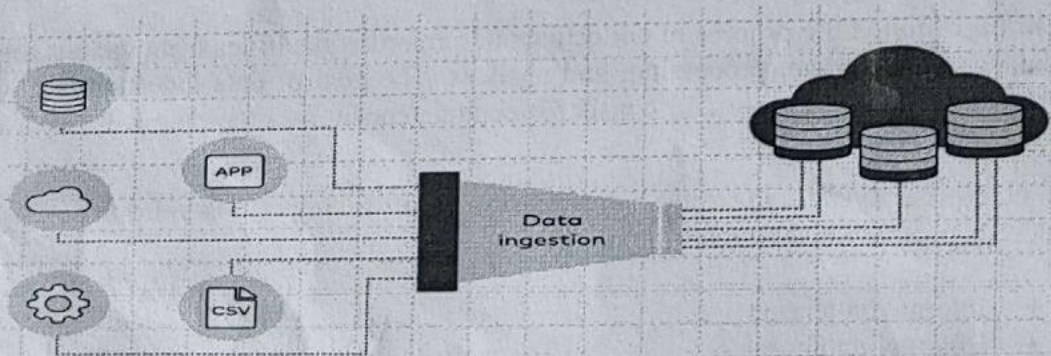
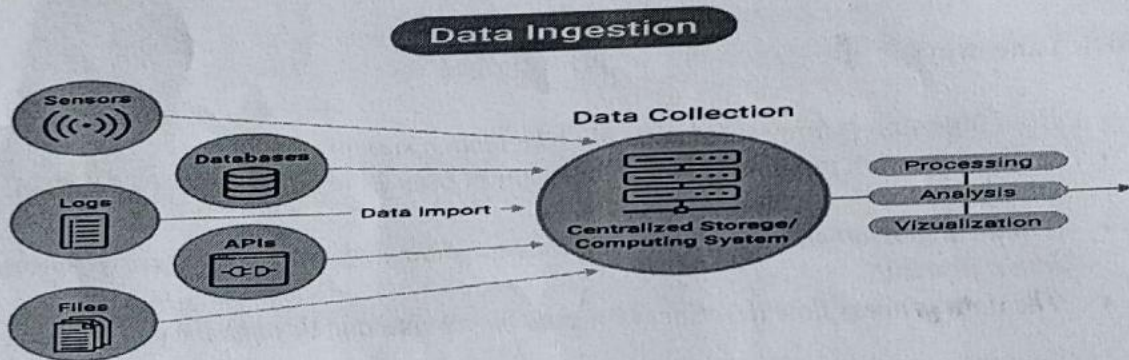
*Customer Relationship Management*

### Types of Data Ingestion:

1. **Batch Ingestion:** Transfers large volumes of data at scheduled intervals. It is cost-effective and well-suited for periodic updates.
2. **Real-Time Streaming:** Continuously streams data into the system, enabling immediate access and analysis. This is critical for applications like IoT monitoring and live dashboards.

Feature	Batch Ingestion	Real-Time Ingestion
Timeliness	Scheduled updates	Immediate updates
Resource Use	Lower computational need	Higher computational need
Use Cases	Historical reporting	Live dashboards, IoT

Data ingestion can occur in batches or in real-time streams. Batch ingestion involves transferring large chunks of data at regular intervals. With streaming ingestion, data is continuously transferred into the system. Typically, real-time streaming ingestion delivers more timely data into the system faster than batch ingestion does.

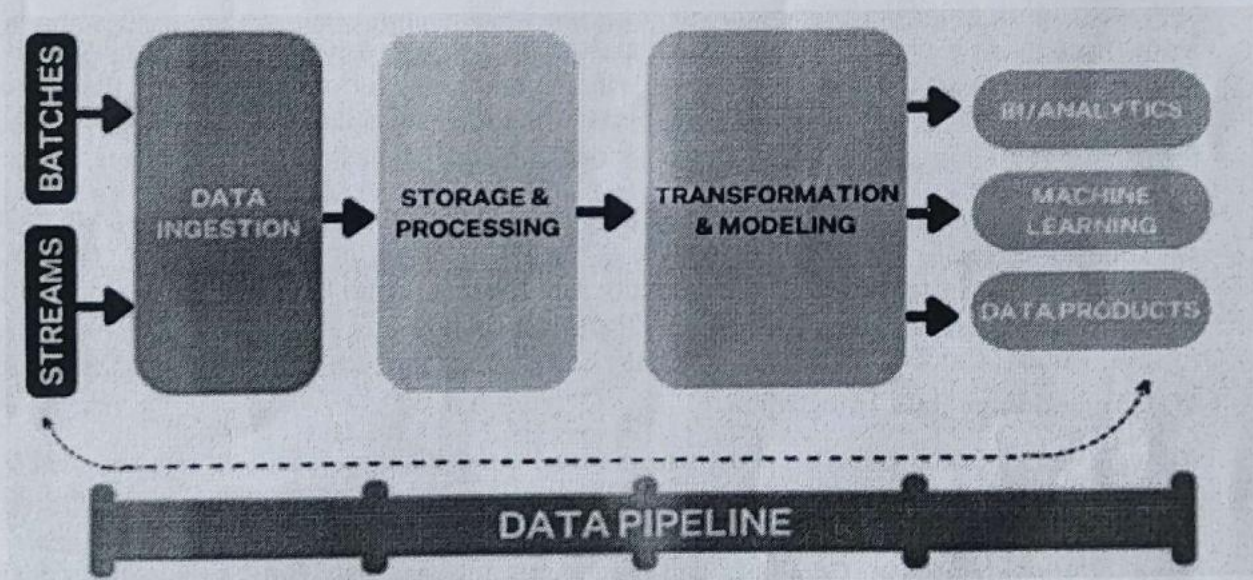
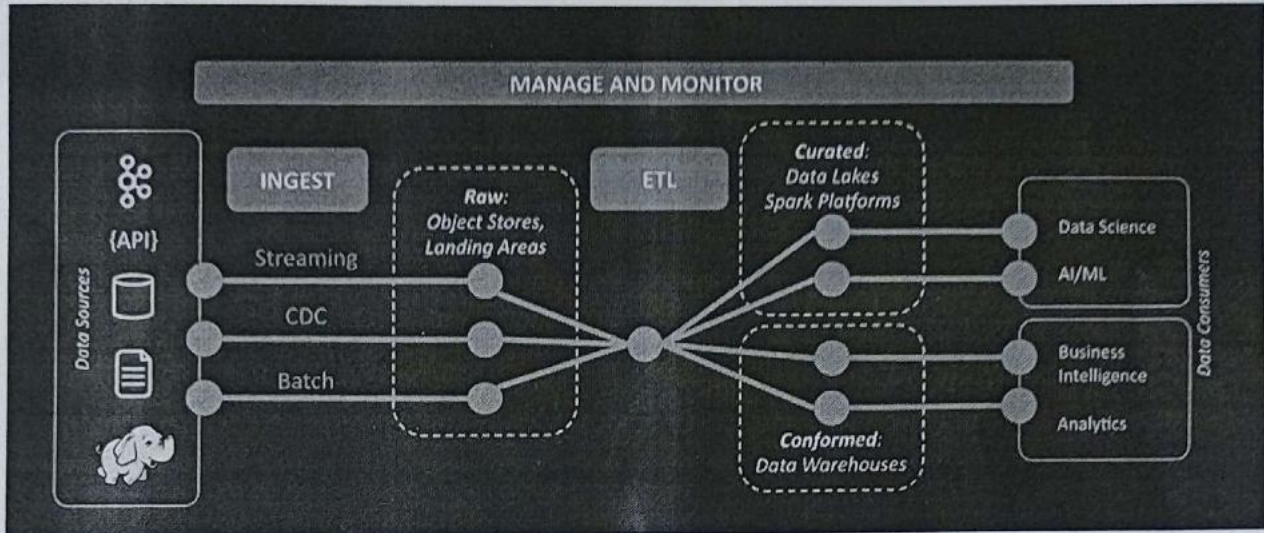


### Data Ingestion Pipelines: Understanding Their Role and Functionality

A data ingestion pipeline is a system that connects multiple data sources to centralized storage, such as a data warehouse or lake. It ensures the smooth movement of data, enabling organizations to structure and organize it for analysis.

## Designing a Robust Data Ingestion Pipeline Architecture

A well-designed **data ingestion pipeline architecture** ensures that data flows seamlessly from its sources to its destination. It must handle diverse formats, support scalability, and incorporate monitoring to maintain data quality.



### Step-by-Step: How Data Ingestion Process Flows Work

#### Process Flow:

1. **Ingestion:** Data is collected from various sources.
2. **Processing:** Data is cleaned, validated, and transformed.
3. **Storage:** Data is centralized in repositories for further analysis.
4. **Analysis:** Data is accessed for insights, reporting, or machine learning applications.

A typical data pipeline has six key data ingestion layers

- **Data ingestion:** This layer accommodates either batched or streamed data from multiple sources.
- **Data storage and processing:** Here, the data is processed to determine the best destination for various analytics, and then stored in a centralized data lake or warehouse.
- **Data transformation and modeling:** Given that ingested data comes in diverse sizes and shapes, not all of it is formally structured (with IDC estimating that 80% of all data is unstructured), this layer transforms all data into a standard format for usability.
- **Data analysis:** This final layer is where users access the data to generate reports and analyses.

Organizations with different data needs may design their data pipelines differently. For instance, a company that only uses batch data may have a simpler ingestion layer. Similarly, a firm that ingests all data in a common format might not need the transformation layer. The data pipeline should be customized to the needs of each organization.

### Data Ingestion vs ETL

Data ingestion and ETL (Extract, Transform, Load) are similar processes but serve different purposes.

**Data ingestion** is a broad term that encompasses the methods used to source, collect, and prepare data for storage or use. It involves bringing data from various sources and making it ready for applications that require it to meet specific formats or quality standards. Typically, the data sources in ingestion processes are not tightly linked to the destination systems.

**ETL**, on the other hand, is a more specific process used in preparing data for storage in data warehouses or data lakes. ETL involves extracting data from one or more sources, transforming it to meet business requirements, and loading it into a destination system. The goal of ETL is often to enable business intelligence, reporting, and analytics.

### Data Ingestion vs Data Integration:

Though related, data ingestion and data integration have distinct roles in a data strategy.

- **Data Ingestion:** Focuses on importing raw data into a centralized system.
- **Data Integration:** Combines and harmonizes ingested data from various sources to create actionable datasets.

Feature	Data Ingestion	Data Integration
Purpose	Importing raw data	Harmonizing and preparing data
Scope	Initial phase	Downstream process

Feature	Data Ingestion	Data Integration
Output	Raw, unstructured data	Structured, usable datasets

**Understanding Data Ingestion Frameworks**

A data ingestion framework outlines the process of transferring data from its original sources to data storage. The right framework enables a system to collect and integrate data from a variety of data sources while supporting diverse data transport protocols.

As noted, data can be ingested in batches or streamed in real time, each approach requiring a unique ingestion framework. Batch data ingestion, a time-honored way of handling large amounts of data from external sources, often involves receiving data in batches from third parties. In other instances, real-time data is accumulated to be ingested in larger batches. A batch data ingestion framework is often less costly and uses fewer computing resources than a streaming framework. However, it's slower and doesn't provide real-time access to the most current data.

In contrast, real-time data ingestion streams all incoming data directly into the data pipeline. This enables immediate access to the latest data but requires more computing resources to monitor, clean, and transform the data in real time. It's particularly useful for data constantly flowing from IoT devices and social media.

Organizations can either design their own data ingestion framework or employ third-party data ingestion tools. Some data ingestion tools support both batch and streamed ingestion within a single framework.

**Understanding Data Ingestion Process Flows**

The data ingestion process flow describes exactly how data is ingested into and flows through a data pipeline. Think of the process flow as a roadmap outlining that data's journey through the system.

When designing a data pipeline, you need to visualize the process flow in advance. This foresight allows the pipeline to be built optimally to handle the anticipated data and its likely usage. Building a pipeline without adequately assessing the process flow could result in an inefficient system prone to errors.

A typical process flow starts at the pipeline's entry point, where data from multiple sources is ingested. The flow continues through layers of the pipeline as the data is stored, processed, transformed, and then analyzed.

**Why High-Quality Data Matters in Data Ingestion?**

Throughout the data pipeline and the process flow, constant monitoring is necessary to ensure the data is clean, accurate, and free from errors. To be useful, data must be:

- Accurate
- Complete
- Consistent
- Timely

- Unique
- Valid

Some experts estimate that 20% of all data is bad—and organizations cannot function with poor-quality, unreliable data. So, any data ingestion process must include robust data quality monitoring, often using third-party tools, to identify poor-quality data and either clean or remove it from the system.

Advanced data pipeline monitoring tools, such as DataBuck from FirstEigen, use artificial intelligence (AI) and machine language (ML) technology to:

- Detect any errors in data ingested into the system
- Detect any data errors introduced by the system
- Alert staff of data errors
- Isolate or clean bad data
- Generate reports on data quality

High-quality data helps an organization make better operational and strategic decisions. If the data is of low quality, business decisions may be compromised.

### 1. What are the steps of a data ingestion pipeline?

A data ingestion pipeline includes these key steps:

- **Data Collection:** Gathering data from multiple sources such as databases, SaaS applications, or IoT devices.
- **Data Transfer:** Moving data to a central system, such as a data lake or warehouse, using batch or real-time ingestion methods.
- **Data Validation:** Ensuring the data is accurate, complete, and meets required standards.
- **Data Storage:** Saving the processed data in a repository for analysis or reporting.

### 2. What is the data ingestion flow?

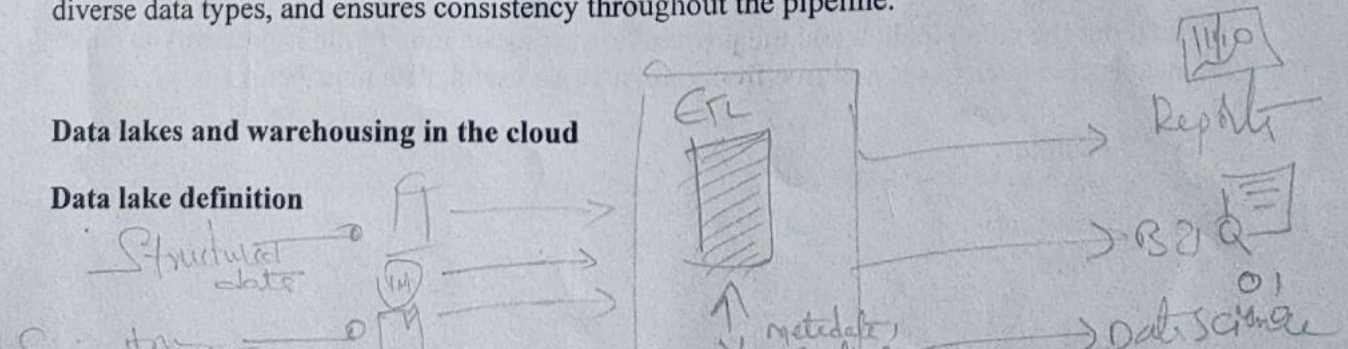
The data ingestion flow outlines how data is collected, validated, and transferred to a centralized system. It typically starts with data sourcing, followed by ingestion into a pipeline, validation for quality, and storage in a data lake or warehouse for further use.

### 3. What is a data ingestion framework?

A data ingestion framework defines the architecture, tools, and processes used to move data from its sources into a central repository. It supports batch and real-time ingestion, integrates diverse data types, and ensures consistency throughout the pipeline.

#### Data lakes and warehousing in the cloud

##### Data lake definition



A data lake is a centralized repository that ingests and stores large volumes of data in its original form. The data can then be processed and used as a basis for a variety of analytic needs. Due to its open, scalable architecture, a data lake can accommodate all types of data from any source, from structured (database tables, Excel sheets) to semi-structured (XML files, webpages) to unstructured (images, audio files, tweets), all without sacrificing fidelity. The data files are typically stored in staged zones—raw, cleansed, and curated—so that different types of users may use the data in its various forms to meet their needs. Data lakes provide core data consistency across a variety of applications, powering big data analytics, machine learning, predictive analytics, and other forms of intelligent action.

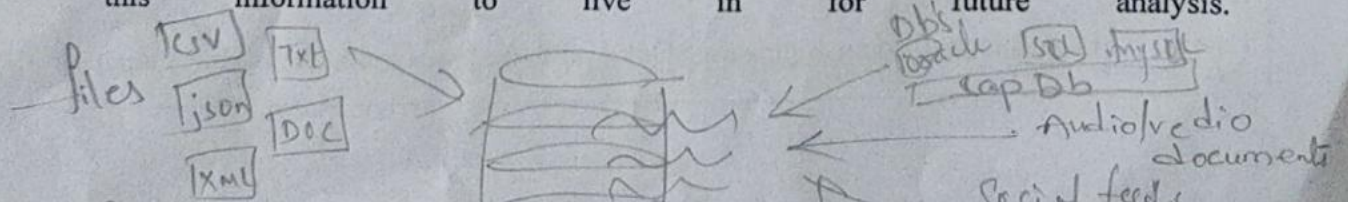
### Why are data lakes important for businesses?

Today's highly connected, insights-driven world would not be possible without the advent of data lake solutions. That's because organizations rely on comprehensive data lakes platforms, such as Azure Data Lake, to keep raw data consolidated, integrated, secure, and accessible. Scalable storage tools like Azure Data Lake Storage can hold and protect data in one central place, eliminating silos at an optimal cost. This lays the foundation for users to perform a wide variety of workload categories, such as big data processing, SQL queries, text mining, streaming analytics, and machine learning. The data can then be used to feed upstream data visualization and ad-hoc reporting needs. A modern, end-to-end data platform like Azure Synapse Analytics addresses the complete needs of a big data architecture centered around the data lake.

### Data lake use cases

With a well-architected solution, the potential for innovation is endless. Here are just a few examples of how organizations across a range of industries use data lake platforms to optimize their growth:

- **Streaming media.** Subscription-based streaming companies collect and process insights on customer behavior, which they may use to improve their recommendation algorithm.
- **Finance.** Investment firms use the most up-to-date market data, which is collected and stored in real time, to efficiently manage portfolio risks.
- **Healthcare.** Healthcare organizations rely on big data to improve the quality of care for patients. Hospitals use vast amounts of historical data to streamline patient pathways, resulting in better outcomes and reduced cost of care.
- **Omnichannel retailer.** Retailers use data lakes to capture and consolidate data that's coming in from multiple touchpoints, including mobile, social, chat, word-of-mouth, and \_\_\_\_\_ in \_\_\_\_\_ person.
- **IoT.** Hardware sensors generate enormous amounts of semi-structured to unstructured data on the surrounding physical world. Data lakes provide a central repository for this information to live in for future analysis.



- **Digital supply chain.** Data lakes help manufacturers consolidate disparate warehousing data, including EDI systems, XML, and JSONs.
- **Sales.** Data scientists and sales engineers often build predictive models to help determine customer behavior and reduce overall churn.

**Data lake vs. data warehouse**

While data lakes and data warehouses are similar in that they both store and process data, each have their own specialties, and therefore their own use cases. That's why it's common for an enterprise-level organization to include a data lake and a data warehouse in their analytics ecosystem. Both repositories work together to form a secure, end-to-end system for storage, processing, and faster time to insight.

A data lake captures both relational and non-relational data from a variety of sources—business applications, mobile apps, IoT devices, social media, or streaming—without having to define the structure or schema of the data until it is read. Schema-on-read ensures that any type of data can be stored in its raw form. As a result, data lakes can hold a wide variety of data types, from structured to semi-structured to unstructured, at any scale. Their flexible and scalable nature make them essential for performing complex forms of data analysis using different types of compute processing tools like Apache Spark or Azure Machine Learning.

By contrast, a data warehouse is relational in nature. The structure or schema is modeled or predefined by business and product requirements that are curated, conformed, and optimized for SQL query operations. While a data lake holds data of all structure types, including raw and unprocessed data, a data warehouse stores data that has been treated and transformed with a specific purpose in mind, which can then be used to source analytic or operational reporting. This makes data warehouses ideal for producing more standardized forms of BI analysis, or for serving a business use case that has already been defined.

	<b>Data lake</b>	<b>Data warehouse</b>
<b>Type</b>	Structured, semi-structured, unstructured	Structured
	Relational, non-relational	Relational
<b>Schema</b>	Schema on read	Schema on write
<b>Format</b>	Raw, unfiltered	Processed, vetted
<b>Sources</b>	Big data, IoT, social media, streaming data	Application, business, transactional data, batch reporting
<b>Scalability</b>	Easy to scale at a low cost	Difficult and expensive to scale
<b>Users</b>	Data scientists, data engineers	Data warehouse professionals, business analysts

Use cases Machine learning, predictive analytics, real-time analytics	Core reporting, BI
---	--------------------

(Or)

**Data Lake:**

Data Lake is the **concept where all sorts of data can be landed at a low cost but exceedingly adaptable storage/zone** to be examined afterward for potential insights. It is another advancement of what ETL/DWH pros called the Landing Zone of data. Only presently we are looking at ALL sorts of information independent of construction, structure, metadata, etc. One of the thoughts behind Data Lake is that presently innovation has made it conceivable to store ALL information that a firm generates/buys (prior it would be a case where the firm HAD to select the pertinent information and store it in a structured distribution center).

**Data Warehouse**

Data Warehouse is essentially a **social database facilitated on cloud or an endeavor centralized computer server**. It collects information from shifted, heterogeneous sources for the most reason for supporting the investigation and choice-making preparation of administration of any business.

A data warehouse is characterized as Subject-oriented, coordinates, time-variant, and non-unstable collection of information in arrange to supply business insights and help within the choice-making process.

**Difference between Data Lake and Data Warehouse**

Data Lake	Data Warehouse
Data is kept in its raw frame in Data Lake and here all the data are kept independent of the source of the information. They are as it was changed into other shapes at whatever point required.	Data Warehouse is composed of data that are extricated from value-based and other measurement frameworks. Here the information isn't in raw shape and is continuously changed and clean.
The most target for Data Lake is Data Researchers, Big Data Engineers, and <u>Machine Learning</u> Engineers who ought to do to profound investigation to form models for commerce such as predictive modeling.	The primary target of Data Warehouse is the operational clients as this information is in organized organize and can give prepared to construct reports. So they are generally utilized for trade intelligence.
The most inputs to data Lake are all sorts of	The primary inputs to Data warehouse

<b>Data Lake</b>	<b>Data Warehouse</b>
information such as organized, semi-structured, and unstructured information. This information dwells in data Lake in their unique form.	are organized information that is coming from value-based and measurements frameworks which are at that point organized within the shape of schemas.
Comprises of raw data that will or might not be curated.	It comprises of curated data which is centralized and is prepared to be used for commerce insights and analytics purpose.
Data is not in normalized form.  The advances that are utilized in data lakes such as <u>Hadoop</u> , Machine Learning are moderately modern as compared to the information warehouse.	Denormalized schemas  Here the technology that's utilized for a data warehouse is older.
A data lake can have all sorts of information and can be utilized with keeping past, show and prospects in mind.	Data Warehouse is concerned, here most of the time is went through on analyzing different sources of the data.
Data in interior of the data lake are profoundly open and can be rapidly updated.	Data in interior of the data warehouse are more complicated and it requires more fetched to bring any changes to them, availability is additionally confined as it were authorized users.

**What is cloud cost optimization?**

Moving data and applications from traditional on-premises data centers to cloud infrastructure offers companies the potential for significant cost savings through accelerating innovation, keeping a competitive edge and better interacting with customers and employees. What's more, IT infrastructure becomes a pay-as-you-go operational expense with most public cloud providers. You can scale your cloud resources up or down to meet demand, and costs will follow. However, cloud services costs can be higher than anticipated, so monitoring and optimizing your cloud spend is critical.

Cloud cost optimization combines strategies, techniques, best practices and tools to help reduce cloud costs, find the most cost-effective way to run your applications in the cloud environment, and maximize business value.

It can be hard to monitor metrics and compare data when using multiple cloud vendors with different dashboards, and overspending can be easy. Whether you use IBM Cloud, Amazon AWS, Google Cloud, Microsoft Azure or some combination of platforms, it's essential to understand, evaluate and optimize what you spend on cloud operations.

### Why do you need cloud cost optimization?

Organizations waste about 32% of their spending on cloud services—a significant sum whether you're a small business or one that spends six or seven digits on the cloud annually. Cloud optimization helps reduce waste and avoid overspending by identifying unused resources and neglected tools.

It's not only about getting costs down. It's also about making sure your costs align with your business goals. In other words, paying more may make sense if you earn more revenue or see more productive activities and profitability from a particular cloud service.

Cloud cost optimization means knowing what your cloud operations cost and making intelligent adjustments so you can control cloud costs without compromising performance.

### Tools for cloud cost optimization

Available cloud cost management tools can help you track bills, features and other configurations, enabling you to optimize costs. Cloud providers offer some tools, including Azure cost management, Google Cloud cost management and AWS cloud financial management tools.

There are also cloud cost tools from independent companies that assess other multiple vendors. For example, IBM® Turbonomic® automates critical actions in real-time, without human oversight, to help you most efficiently use compute, storage and network resources. These tools can work across multiple clouds and create reports showing the combined, multicloud data.

OR

Cloud cost optimization for AI involves strategies like **right-sizing** compute (GPUs/CPU), using **Spot/Reserved instances**, implementing data lifecycle management for cheaper storage tiers (Archive/Glacier), automating shutdowns for idle resources, and leveraging AI/ML for workload scheduling and anomaly detection to maximize performance and minimize waste. Key areas are compute (spot instances, rightsizing), storage (tiering, deletion), and workload management (scheduling, segmentation).

### Compute Optimization

- **Right-Sizing & Autoscaling:** Match instance types (especially GPUs) to actual needs; scale down or off during idle periods.

- **Pricing Models:** Use Spot Instances for fault-tolerant training/inference; leverage **Reserved Instances (RIs)** or **Savings Plans** for predictable workloads.
- **Serverless & Containers:** Use serverless functions for sporadic tasks or container orchestration (Kubernetes) for efficient resource pooling.
- **Efficient Architectures:** Choose cost-effective AI services (e.g., managed endpoints vs. self-hosted).

### Storage Optimization

- **Data Lifecycle Management (DLM):** Automatically move data from expensive hot storage (S3 Standard) to cheaper tiers (Infrequent Access, Archive, Deep Archive) as it ages.
- **Data Tiering & Archiving:** Archive old datasets, models, and logs to Glacier/Archive tiers.
- **Clean Up:** Regularly delete unused volumes, snapshots, and temporary files.
- **Intelligent Tiering:** Use features like S3 Intelligent-Tiering to automatically move data based on access patterns.

### Workload & Data Management

- **Workload Classification:** Separate experimental, training, and inference workloads for tailored resource allocation.
- **Scheduling:** Schedule non-urgent jobs (e.g., large training runs) during off-peak hours or when Spot prices are low.
- **Caching & Segmentation:** Cache results for inference; break large workflows into smaller, schedule-able segments.
- **Resource Tagging:** Tag resources to track ownership, usage, and costs, making it easier to identify waste.

### Monitoring & Governance

- **FinOps:** Implement FinOps practices with continuous monitoring and cost governance policies.
- **Automated Alerts:** Set up alerts for cost anomalies and idle resources.
- **Build vs. Buy:** Evaluate using managed AI APIs versus building custom models to balance cost and control.

- ⇒ Amazon Simple Storage Service (S3) is an online cloud based data storage infrastructure for storing & retrieving any amount of data.
  - ⇒ S3 provides highly available, scalable, fast, fully redundant & affordable storage infrastructure.
  - ⇒ The data stored on S3 is organized in the form of buckets.
  - ⇒ You must create a bucket before you can store data on S3.
  - ⇒ S3 console provides simple wizards for creating a new bucket & uploading files.
  - ⇒ You can specify the redundancy & encryption options & access permissions.
- root folders.  
 ↓  
 Inside sub folder we call objects.
- ⇒ bucket name should be "unique".

### Windows Azure Storage

Windows Azure Storage is the cloud storage services from Microsoft.

- ⇒ It provides various storage services such as blob storage service, table service & queue service.
  - ⇒ The blob storage services allow storing unstructured binary data large objects (blobs).
  - ⇒ blobs are organized into containers two kinds of blobs can be stored.
    - block blobs
    - page blobs
  - ⇒ A blob can be subdivided into some number of blocks if a failure occurs while transferring a block blob, retransmission can resume with the most recent block rather than sending the entire blob again.
- ... divided into number of pages and are designed

BigQuery :- BigQuery is Google cloud's fully managed, serverless <sup>Engine for</sup> Enterprise data warehouse, which decouples Compute & Storage. <sup>It leverages a powerful SQL Query</sup>  
- large-scale data analysis.

Google BigQuery is a serverless data warehouse used for data analysis.

Key Char

- 1) Service Type :- Serverless data warehouse.
  - 2) Data structure :- Optimized for structured & semi-structured data
  - 3) Integration :- Query data can be stored directly in GCS. <sup>its BigQuery omni</sup>
  - 4) Functionality :- BigQuery separates <sup>its</sup> Compute Engine from <sup>its</sup> storage. <sup>without provisioning infrastructure.</sup>
- Use cases :- BI, ML, Real-time analytics & running complex SQL queries over massive datasets.

# ~~Elastic~~ Computing (EC<sub>2</sub>)

Elastic Computing is a concept closely related to utility computing & cloud computing...

It refers to the ability to dynamically adjust the amount of computing resources allocated to a workload based on changing demands.

Elastic Computing is the ability of CSP to swiftly scale the usage of resources such as storage, infrastructure, memory, RAM, CPU, <sup>IO</sup> bandwidth, ~~etc~~ etc.

To <sup>→ upgrading the # of resources</sup> scale up & down <sup>→ addition of more resources of same type</sup> to adapt the changing resource demands & dynamically meet workload requirements.

⇒ Elastic Computing is a part of CC that entails dynamically managing the cloud server.

In elastic computing environments, resources such as processing power, memory & storage can be automatically scaled up (↑) down in response to fluctuations in workload (↓) user demand.

This elasticity allows organizations to efficiently use resources ensuring that they have enough capacity to handle peak loads without overprovisioning & wasting resources during periods of lower demand. <sup>↓</sup> leads to higher capital expenditure than <sup>required</sup>.

Elasticity is a key feature of CC platforms, where resources can be <sup>under</sup> provisioned & <sup>stop</sup> de-provisioned automatically through tools like auto-scaling policies. This flexibility enables business to optimize costs, improve performance & maintain

Big Query is Google cloud's fully managed serverless data warehouse. It leverages a powerful SQL engine.

Variation in workloads leads to rapid, traditional approaches are unable to keep track with the demand which leads to either over-provisioning (wasting resources) or under-provisioning (causing performance issues).

BigQuery is a serverless data warehouse that allows you to store and analyze massive amounts of data. It is designed for high performance and scalability. BigQuery uses a distributed architecture to process queries across many nodes. This allows it to handle large datasets and complex queries efficiently. BigQuery also offers a variety of features, including data partitioning, clustering, and caching, to optimize query performance. Additionally, BigQuery provides a rich set of APIs and integrations with other Google Cloud services, making it easy to integrate into existing data pipelines and workflows.

Virtualization :- is a technology in cloud computing.

Def: Virtualization refers to the partitioning the resources of a physical system into multiple virtual resources.  
(Computing, Storage, N/w & Memory)

⇒ It is a key enabling technology of cloud computing & allows pooling of resources.

⇒ In CC resources are pooled to serve multiple users

Using "Multi-tenancy".

⇒ Multi-tenant aspects of the cloud allow multiple users to be served by the "same physical hardware".

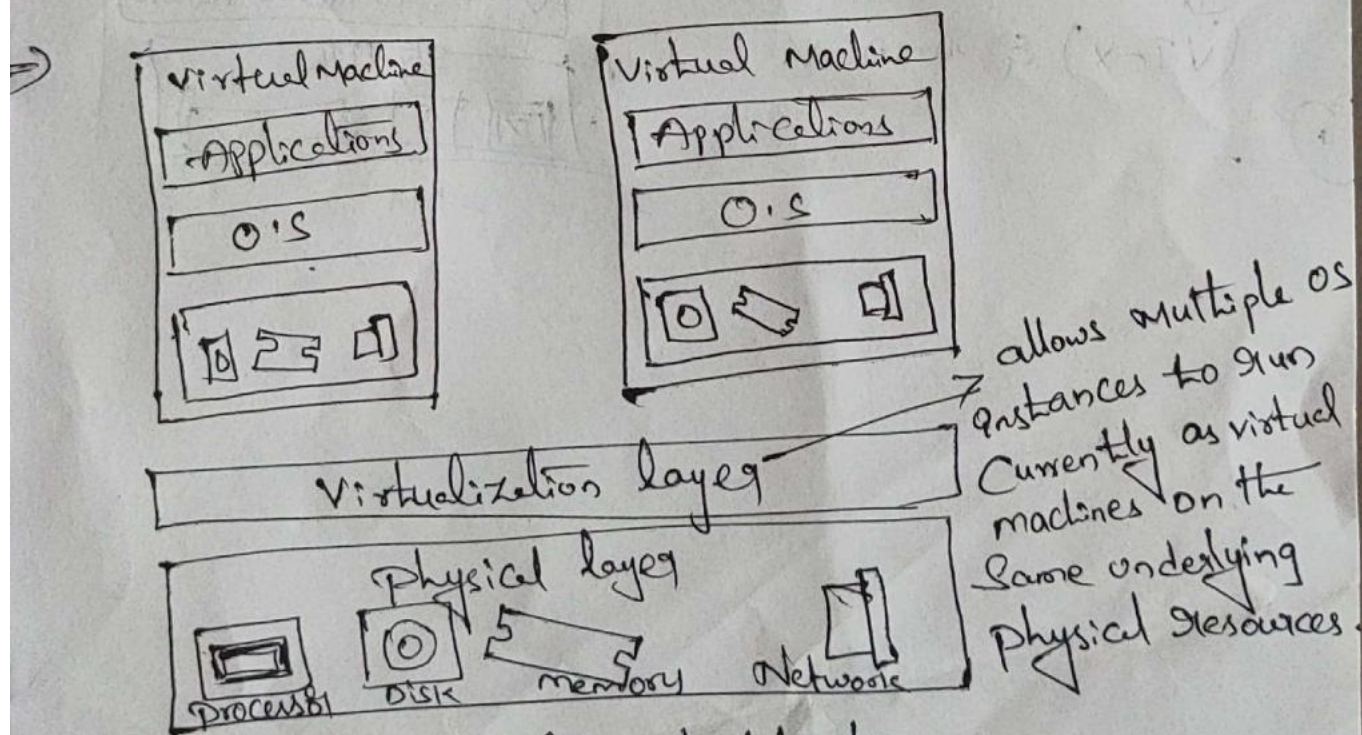


Fig:- Virtualization Architecture.

⇒ The physical resources such as computing, storage, memory & N/w resources are virtualized.

⇒ The virtualization partitions the physical resources into multiple virtual machines.

# Hypervisor:-

Virtualization layer consists of a hypervisor (or) VMM (Virtual Machine Monitor).

The hypervisor presents a virtual platform to a Guest OS.

↓  
that is installed in a virtual machine in addition to the host (or) main OS.

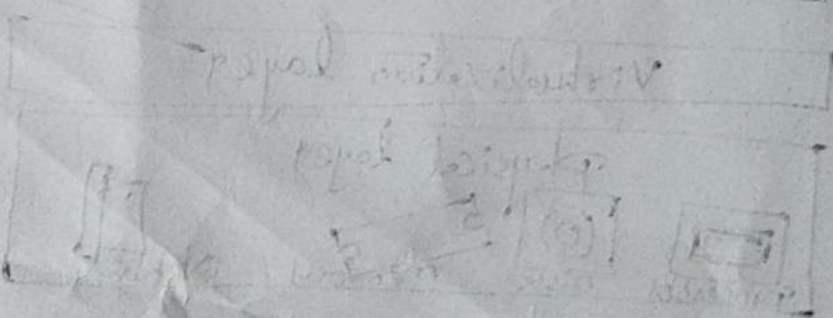
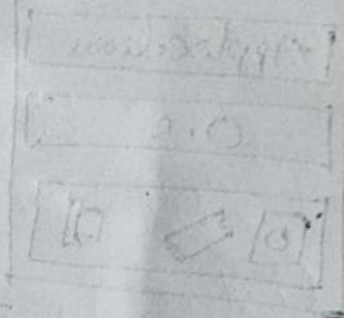
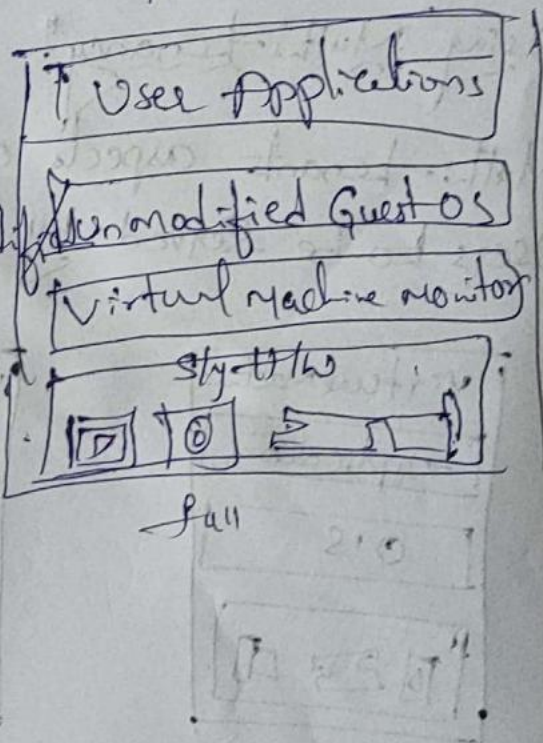
Various forms of virtualization approaches exists:

1) Full virtualization →

2) Para " "

3) H/W " "

↓  
(VT-X) & AMD-V



## UNIT-III

### ML Services on AWS – Amazon SageMaker

**Amazon SageMaker** is a fully managed **Machine Learning (ML) service** provided by AWS that enables developers and data scientists to **build, train, deploy, and manage ML models at scale** quickly and efficiently.

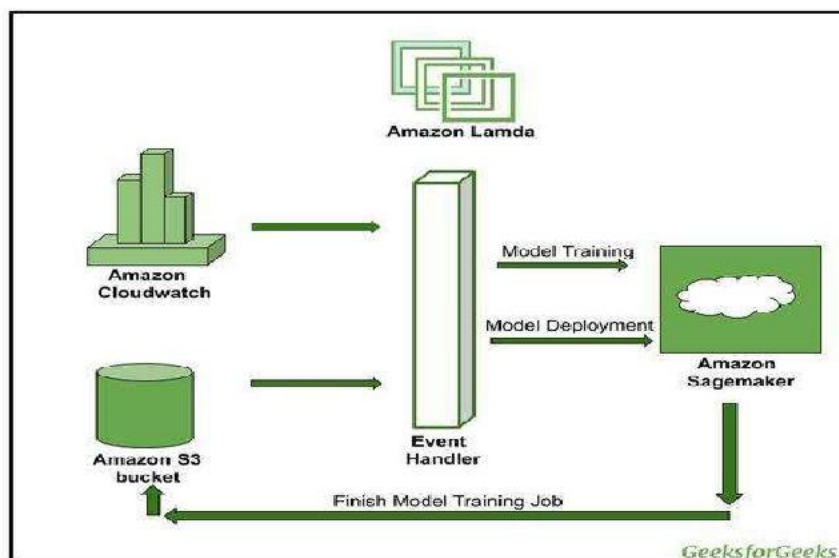
Unlike traditional ML development, which requires managing complex infrastructure for training and hosting, SageMaker abstracts away the heavy lifting. It provides a unified toolset for every step of the ML lifecycle, from labeling raw data to monitoring production models for drift.

### Machine learning in AWS SageMaker

Machine learning (ML) within AWS SageMaker follows a cyclical process that requires both workflow management tools and specialized hardware to handle large data sets. Typically, ML models are developed in two main stages: training and inference.

In the training phase, the system learns to identify patterns in the data, allowing it to predict outcomes based on similar patterns in the future. After training, the model moves to inference, where it analyzes new data to make predictions. Once data scientists have fine-tuned the model, development teams then transform the trained model into application program interfaces (APIs) that can be integrated into products or services.

Many organizations face challenges in AI development due to the costs of hiring experts and maintaining the necessary infrastructure. AWS SageMaker addresses these challenges by offering integrated tools that automate manual tasks, reduce human error, and minimize hardware expenses. The platform provides a suite of ML modeling tools within an easy-to-use framework. With SageMaker templates, businesses can quickly build, train, host, and deploy machine learning models at scale in the AWS cloud.



SageMaker is built around three distinct stages of the ML lifecycle. Importantly, **you can use these independently**. You can train a model in SageMaker and deploy it elsewhere, or train a model on your laptop and deploy it to SageMaker.

### 1. Build (SageMaker Studio & Notebooks)

- **SageMaker Studio:** A fully integrated development environment (IDE) for ML. It provides a single web-based interface for all ML steps.
- **Notebook Instances:** Fully managed EC2 instances running JupyterLab. You can spin them up in minutes to explore data and write code.

### 2. Train (Training Jobs)

- **On-Demand Infrastructure:** When you start a "Training Job," SageMaker spins up a cluster of EC2 instances (e.g., highly powerful GPU instances like p3.16xlarge), loads your data from S3, runs your training script, saves the model artifacts to S3, and then automatically terminates the cluster.
- **Cost Efficiency:** You only pay for the exact seconds the training cluster is running.

### 3. Deploy (Inference Endpoints)

- **Model Hosting:** SageMaker takes your trained model artifacts from S3 and deploys them to a fleet of instances behind a secure HTTP endpoint. This allows your applications to request predictions in real-time.

## AWS SageMaker Workflow

1. **Data Preparation:** The first step in the workflow is to prepare the data for training the machine learning model. This includes tasks such as collecting, cleaning, and transforming data into the appropriate format.
2. **Model Building:** Once the data is prepared, the next step is to build the machine learning model. SageMaker provides a variety of pre-built algorithms and frameworks, or users can bring their own custom algorithms.
3. **Model Training:** After the model is built, the next step is to train it using the prepared data. SageMaker provides a range of options for training, including distributed training on multiple instances for faster results.
4. **Model Optimization:** Once the model is trained, the next step is to optimize it for performance. This includes tasks such as fine-tuning hyperparameters and optimizing the model's architecture.
5. **Model Deployment:** Once the model is optimized, the next step is to deploy it for use in a production environment. SageMaker provides options for deploying models to various endpoints, including [Amazon EC2 instances](#), [Lambda functions](#), and [API Gateway](#).
6. **Model Monitoring:** Once the model is deployed, the next step is to monitor its performance in real time. SageMaker provides built-in monitoring tools that track the model's performance metrics and detect anomalies.
7. **Model Management:** Finally, once the model is in production, it's important to manage it over time. This includes tasks such as updating the model with new data, retraining the model periodically, and ensuring that it remains performant over time.

## Key Features of Amazon SageMaker

### 1. Fully Managed Environment

- No need to manage servers or infrastructure
- Handles provisioning, scaling, and maintenance automatically

2. **Integrated Development Environment**
  - Provides **SageMaker Studio** and **Jupyter Notebooks**
  - Supports Python-based ML frameworks
3. **Built-in Algorithms**
  - Pre-built, optimized algorithms for:
    - Linear Regression
    - XGBoost
    - K-Means
    - Image Classification
    - Text Classification
4. **Framework Support**
  - Supports popular ML frameworks:
    - TensorFlow
    - PyTorch
    - MXNet
    - Scikit-learn

### Advantages of SageMaker

- Faster ML model development
- Reduced operational complexity
- Scalable and cost-effective
- High availability and reliability

### Use Cases:

Organizations across various industries use SageMaker to address a wide range of challenges:

- **Fraud Detection:** Real-time scoring of transactions in financial institutions.
- **Personalization:** Building recommendation engines for e-commerce or media streaming services.
- **Predictive analytics:** Analyzing sensor data to predict equipment failures before they happen.
- **Generative AI:** Building and fine-tuning custom generative AI applications and large language models (LLMs).
- **Image and speech recognition**
- **Natural Language Processing (NLP)**
- **Recommendation systems**

**The next generation of Amazon SageMaker consists of two primary components:**

1. **Amazon SageMaker Unified Studio**, which provides an integrated experience to use all your data and tools for analytics and AI

2. **Data and AI governance**, which applies enterprise-level security and data management with built-in governance throughout the entire data and AI lifecycle

Additionally, SageMaker is built upon an open lakehouse architecture that unifies access to all your data across Amazon Simple Storage Service ([Amazon S3](#)) data lakes, [Amazon Redshift](#) data warehouses, and other external sources



### ML Services on AZURE ML:

**Azure Machine Learning (Azure ML)** is a comprehensive, enterprise-grade, cloud-based platform provided by [Microsoft](#) that accelerates the end-to-end machine learning and deep learning project lifecycle. It provides a managed environment for building, training, deploying, and managing models at scale, supporting both low-code/no-code and code-first development approaches.

### Key Features and Components

Azure ML is designed for individuals and teams looking to implement Machine Learning Operations (MLOps) and integrate ML models into secure production environments.

- **Diverse Tooling:** Users can work with their preferred tools, including the Azure ML studio web interface, Python SDK (v2), Azure CLI (v2), and REST APIs.
- **Code-first & No-code Options:**
  - **Notebooks:** Write and run custom Python or R code in managed Jupyter Notebook servers.
  - **Automated Machine Learning (AutoML):** Automatically identifies suitable algorithms and tunes hyperparameters, ideal for users without extensive data science knowledge or who want to accelerate the process.

- **Designer:** A drag-and-drop visual interface for building, testing, and deploying ML pipelines without writing code.
- **Open-Source Framework Support:** Azure ML supports popular open-source ML and deep learning frameworks like **PyTorch**, **TensorFlow**, **scikit-learn**, and more, allowing flexibility and avoiding vendor lock-in.
- **Scalable Compute Resources:** Provides on-demand access to a wide range of compute resources, including CPU and high-performance GPU virtual machines and clusters, necessary for compute-intensive deep learning tasks.
- **MLOps and Lifecycle Management:** Offers integrated tools for the entire ML lifecycle, including tracking experiments, data versioning, creating reusable pipelines, monitoring model performance, and implementing CI/CD workflows with GitHub Actions or Azure DevOps.
- **Enterprise Security:** Integrates with Azure security features like Virtual Networks and Azure Key Vault for secure, compliant ML projects.
- **Generative AI:** Includes tools like a model catalog (featuring models from Azure OpenAI Service, Mistral, Meta, etc.) and prompt flow to streamline the development cycle of applications powered by Large Language Models (LLMs).

### Advantages

- Fully managed and scalable platform
- Faster ML development using AutoML
- Seamless integration with Azure services
- Enterprise-grade security

### Applications

Azure ML is used in **predictive analytics, image processing, NLP, fraud detection, and recommendation systems.**

### ML SERVICES ON GCP VERTEX AI:

Google Cloud's Vertex AI is a fully managed, unified platform for the entire machine learning (ML) and deep learning (DL) development lifecycle. It consolidates various AI tools and services into a single environment, enabling users with varying levels of expertise to build, deploy, and scale ML models efficiently.

**Vertex AI** is GCP's unified machine learning platform designed to manage the complete ML lifecycle. It provides tools for **data preparation, model training, evaluation, and deployment** through a single interface. Vertex AI supports **AutoML and custom training** using frameworks like TensorFlow and PyTorch. It offers scalable infrastructure, built-in pipelines, and secure model deployment.

## Key Features and Tools

Vertex AI offers various tools for each stage of the ML workflow:

- **Data Preparation & Storage:** This includes tools for collecting, cleaning, and transforming data. It integrates with services like Cloud Storage and BigQuery.
- **Model Training:**
  - **AutoML:** This is a low-code option that automates model creation for those with limited ML expertise.
  - **Custom Training:** Provides full control for advanced users and supports frameworks like TensorFlow and PyTorch.
  - **Generative AI:** Access to Google's foundation models, including Gemini, in the **Vertex AI Studio** for tuning and deployment.
- **Model Management (MLOps):**
  - **Vertex AI Pipelines:** This orchestrates and automates ML workflows, making them repeatable and scalable.
  - **Vertex AI Model Registry:** Manages the lifecycle of different model versions.
  - **Vertex AI Feature Store:** Serves as a centralized repository for serving, sharing, and reusing ML features.
  - **Vertex AI Model Monitoring:** Tracks deployed models for performance degradation, data skew, and drift.
- **Model Deployment & Inference:** Simplifies deploying models to production for online or batch predictions. It has capabilities like Vertex Explainable AI for understanding model outputs.
- **Development Environments:** **Vertex AI Workbench** offers managed JupyterLab notebooks, and **Colab Enterprise** provides a collaborative, zero-config environment for coding.

## Advantages

- **Unified Platform:** Manages the entire end-to-end ML lifecycle in one place, reducing complexity.
- **Scalability:** Automatically scales computing resources to meet the demands of growing data and complex models.
- **Flexibility:** Accommodates different skill levels, from business users using AutoML to ML experts needing custom training and control.
- **Integration:** Seamlessly integrates with other Google Cloud services (e.g., BigQuery, Cloud Storage) and supports open-source frameworks.

## Advantages

- Unified and user-friendly platform
- Highly scalable and cost-efficient

- Strong integration with GCP services
- High model performance using Google's AI infrastructure

## Applications

Vertex AI is widely used in **image recognition, NLP, demand forecasting, fraud detection, recommendation systems, and predictive analytics.**

## Use Cases

Vertex AI is used across industries for tasks such as:

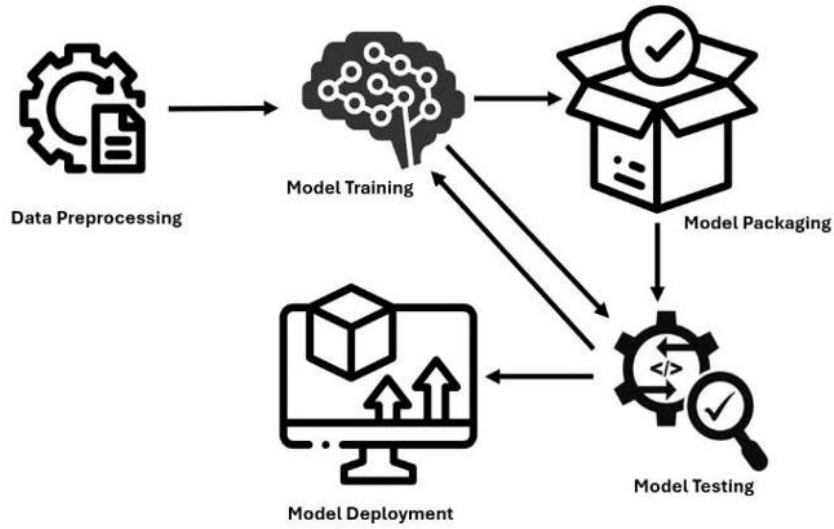
- Predicting customer churn and detecting fraudulent transactions in financial services.
- Demand forecasting and personalized recommendations in retail.
- Optimizing treatments and predicting patient outcomes in healthcare.
- Computer vision applications, like object detection and video analysis.

All three platforms provide **end-to-end ML capabilities**, but:

- **SageMaker** is ideal for large-scale and customizable ML solutions
- **Azure ML** is best for enterprise and Microsoft-centric ecosystems
- **Vertex AI** excels in AutoML and unified ML workflows

## 2. Training and Deploying Models on Cloud:

Model deployment is the process of trained models being integrated into practical applications. This includes defining the necessary environment, specifying how input data is introduced into the model and the output produced, and the capacity to analyze new data and provide relevant predictions or categorizations. Let us explore the process of deploying models in production.



## Step-by-Step Flow Diagram: Training and Deploying Models on Cloud

Start



Data Collection



Data Storage (Cloud Storage / Data Lake)



Data Preprocessing



Feature Engineering



Model Selection

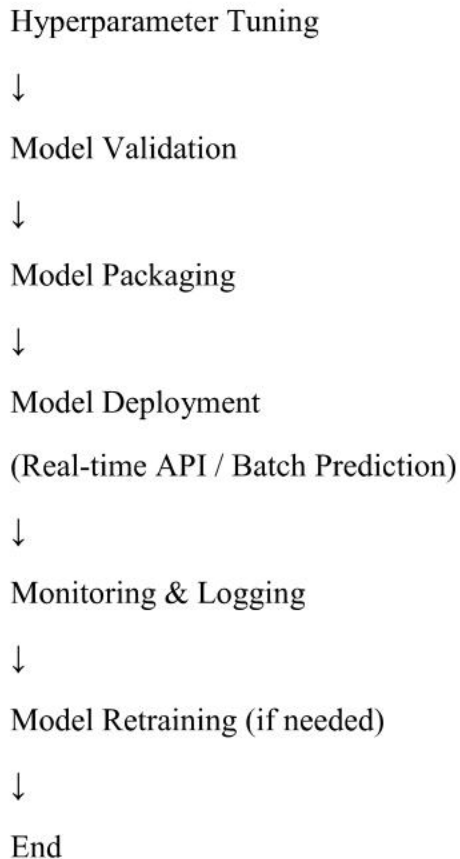


Model Training (Cloud Compute: CPU / GPU / TPU)



Model Evaluation





1. **Choose the Right Cloud Provider:** Select a cloud provider based on your specific needs, such as AWS, Google Cloud, Microsoft Azure, or IBM Cloud. Each provider offers unique features and services for model deployment and management.
2. **Model Training:** Train your model on the cloud using scalable compute instances, such as GPUs, TPUs, and high-memory instances. Cloud platforms provide tools and services to speed up the training process.
3. **Deployment:** Deploy the trained model to the cloud using the chosen provider's services. This can involve creating an endpoint for inference, setting up a SageMaker endpoint, or using Azure Machine Learning for model management.
4. **Monitoring and Maintenance:** Continuously monitor the model's performance and make necessary adjustments. Cloud platforms often provide tools for model monitoring and maintenance, ensuring the model remains effective and up-to-date.
5. **Use Cases:** Deploy your model in various use cases, such as predictive analytics, customer behavior prediction, or real-time decision-making. The cloud's scalability and flexibility make it suitable for a wide range of applications.

By following these steps and leveraging the capabilities of cloud platforms, you can effectively deploy AI and ML models in the cloud, transforming data science projects into real-world applications.

## Google Vertex AI

Google [Vertex AI](#) is a fully managed ML platform that integrates Google Cloud's robust AI tools to simplify building, deploying and scaling ML models. It supports custom trained models, AutoML-generated models and Google's foundation models (e.g., PaLM, Imagen).

### Deployment Features:

- **Model Serving:** Supports scalable online prediction with autoscaling endpoints and batch prediction. Allows deployment on GPUs, CPUs or TPUs.
- **Vertex Pipelines:** Automates end-to-end ML workflows including deployment.
- **Integration:** Deep integration with Google Cloud services such as BigQuery (data), Dataflow (streaming) and [Explainable AI](#) for interpretability.
- **AutoML & Custom Models:** Users can choose AutoML for automated model creation or deploy custom [TensorFlow](#), PyTorch or other models.
- **Multimodal Support:** Native support for text, vision, audio and structured data models.

**Use Cases:** Data-rich enterprises needing complex, customizable workflows and integration with Google's ecosystem, heavy on data analytics and large batch/online predictions.

## Step-by-Step Deployment on Google Vertex AI

- **Initialize Vertex AI Client and Environment**

Set your Google Cloud project and region, initialize the Vertex AI client in Python:



```
from google.cloud import aiplatform
```

### # Initialize Vertex AI

```
aiplatform.init(project='YOUR_PROJECT_ID', location='us-central1')
```

- **Upload/Register Your Model**

Upload your trained model artifacts (e.g., TensorFlow SavedModel or PyTorch model) stored in Google Cloud Storage (GCS) to Vertex AI Model Registry



```
model = aiplatform.Model.upload(  
    display_name='my_custom_model',  
    artifact_uri='gs://your-bucket/path-to-model/',  
    serving_container_image_uri='us-docker.pkg.dev/vertex-ai/prediction/tf2-cpu.2-3:latest'  
)
```

Replace `serving_container_image_uri` with a container matched to your model framework (TensorFlow, PyTorch, AutoML, etc.)

- **Create or Get an Endpoint**

Create a new endpoint to deploy your model, which provides a URL to send prediction requests



```
endpoint = aiplatform.Endpoint.create(display_name='my_model_endpoint')
```

Alternatively, you can use an existing endpoint.

- **Deploy Model to Endpoint**

Bind your model to the endpoint with configured compute resources (machine type, replicas, autoscaling)



```
endpoint.deploy(  
    model=model,  
    deployed_model_display_name='my_model_deployment',  
    machine_type='n1-standard-4', # e.g. CPU machine; GPU/TPU types are available  
    min_replica_count=1,  
    max_replica_count=3,  
    traffic_split={"0": 100},  
)
```

The deployment process provisions the necessary compute resources to serve online predictions.

- **Send Prediction Requests**

Make online prediction calls by sending HTTP POST requests to the endpoint URL (or use the SDK)



```
response = endpoint.predict(instances=[YOUR_INPUT_DATA])  
print(response)
```

## 2. Amazon SageMaker

[AWS SageMaker](#) is a comprehensive ML development and deployment service optimized for a broad range of AI workloads, including classic ML and deep learning models.

### Deployment Features:

- **Multi-Option Deployment:** Supports real-time inference endpoints, batch transform jobs and asynchronous inference for large payloads.
- **Serverless Inference:** Automatically provisions compute resources for prediction workloads without infrastructure management.

- **Multi-Framework Support:** Compatible with TensorFlow, PyTorch, MXNet, Scikit-learn and more, with pre-built containers.
  - **Model Monitoring:** Integrated continuous monitoring detects model drift and data quality issues.
  - **JumpStart:** Provides pre-trained models and solutions with one-click deployment.
  - **Pipeline Automation:** SageMaker Pipelines enable full CI/CD for ML models.
- Use Cases:** Teams requiring customizable, end-to-end ML lifecycle management and infrastructure control with extensive monitoring and tuning.

### Step-by-Step Deployment on AWS SageMaker

- **Upload the Model to S3:** Save your model and upload it to an Amazon S3 bucket.



```
aws s3 cp model.pkl s3://your-bucket-name/
```

- **Create a SageMaker Endpoint:** Use the SageMaker console or SDK to create an endpoint for inference.



```
import boto3
```

```
sagemaker = boto3.client('sagemaker')
response = sagemaker.create_model(
    ModelName='my-ml-model',
    PrimaryContainer={
        'Image': 'your-container-image',
        'ModelDataUrl': 's3://your-bucket-name/model.pkl'
    }
)
```

- **Deploy and Test:** Deploy the model as an endpoint and make predictions via the SageMaker API.

### Microsoft Azure Machine Learning

Azure ML is Microsoft's enterprise ML platform offering build, train and deployment capabilities with integrated MLOps and data governance.

#### Deployment Features:

- **Flexible Deployment:** Supports online endpoints (real-time) and batch inference with autoscaling and GPU/CPU support.
- **AutoML & Custom Models:** Enables automated model creation plus custom model deployment.

- **Integration:** Tight coupling with Azure tools like [Power BI](#), [Azure DevOps](#), [Cognitive Services](#) and the broader Microsoft business ecosystem.
- **Model Registry:** Centralized catalog for versioning and lifecycle management.
- **MLOps Pipelines:** Supports CI/CD and governance automation.

**Use Cases:** Enterprises leveraging Microsoft products with emphasis on governance, compliance and hybrid deployments.

### Step-by-Step Deployment on Azure ML

- **Register the Model:** Use the Azure ML workspace to register your trained model.
- **Create a Scoring Script:** Write a script that defines how the model should handle requests.



```
def init():
    global model
    model = load_model('model.pkl')
```

```
def run(input_data):
    data = json.loads(input_data)
    return model.predict(data['features'])
```

- **Deploy as a Web Service:** Use Azure CLI or SDK to deploy the model.



```
az ml model deploy --model-id my_model --name my-service --runtime python --deployment-type aks
```

- **Test the Deployed Endpoint:** Send HTTP POST requests with appropriately formatted input data (usually JSON) to the endpoint's URL.

## 3. AutoML and Custom ML Model Workflows:

Machine learning workflows define the **step-by-step process** used to build, train, evaluate, deploy, and maintain ML models. Based on the level of automation and user control, ML workflows are commonly classified into **AutoML workflows** and **Custom ML model workflows**.

### 1. AutoML Workflow

**AutoML (Automated Machine Learning)** automates the entire ML pipeline with minimal human intervention. It is designed for users who have limited ML expertise and want **quick and accurate models**.

### AutoML Workflow Steps

1. **Data Collection**
  - User provides structured or unstructured datasets.
  - Data is uploaded to cloud storage.
2. **Data Preprocessing**
  - Automatic data cleaning, normalization, and handling of missing values.
  - Feature scaling and encoding are performed automatically.
3. **Feature Engineering**
  - AutoML automatically generates and selects relevant features.
4. **Model Selection**
  - Multiple algorithms are tested automatically.
  - Best-performing model is selected.
5. **Hyperparameter Tuning**
  - AutoML optimizes hyperparameters to improve performance.
6. **Model Evaluation**
  - Models are evaluated using metrics such as accuracy, precision, recall, or RMSE.
7. **Model Deployment**
  - Best model is deployed as a cloud service or API.
8. **Monitoring and Retraining**
  - Model performance is continuously monitored.
  - Retraining is triggered when performance degrades.

### Advantages of AutoML

- Minimal coding required
- Faster model development
- Suitable for non-experts
- Efficient for standard ML problems

### Limitations of AutoML

- Limited customization
- Less control over model internals
- Not ideal for complex or research-level problems

## 2. Custom ML Model Workflow

**Custom ML workflows** provide full control over every stage of the ML lifecycle. They are used by experienced data scientists for **complex, specialized, or high-performance applications**.

## Custom ML Workflow Steps

1. **Problem Definition**
  - Clearly define objectives, constraints, and evaluation metrics.
2. **Data Collection**
  - Collect data from databases, APIs, sensors, or logs.
3. **Data Preprocessing**
  - Manual data cleaning, normalization, and transformation.
4. **Feature Engineering**
  - Domain-specific feature extraction and selection.
5. **Model Selection**
  - Choose appropriate ML/DL algorithms.
6. **Model Training**
  - Train models using custom code and frameworks.
7. **Hyperparameter Tuning**
  - Manually or programmatically optimize parameters.
8. **Model Evaluation**
  - Evaluate model performance using test datasets.
9. **Model Deployment**
  - Deploy models as APIs or batch jobs.
10. **Monitoring and Retraining**
  - Track accuracy, latency, and drift.
  - Retrain models with new data.

## Advantages of Custom ML

- Full flexibility and control
- Higher performance for complex tasks
- Suitable for deep learning and research

## Challenges of Custom ML

- Requires expert knowledge
- Longer development time
- Higher operational complexity

## 4.Gpu/TPUs FOR MODEL TRAINING:

**Cloud-based GPU platforms provide powerful infrastructure for machine learning, enabling scalable and efficient training of AI models without the need for on-premise hardware.**

### Key Features of Cloud GPU Platforms

1. **Performance and Hardware:** Look for platforms that offer the latest NVIDIA GPUs (such as A100, H100, and H200) which can significantly accelerate deep learning tasks. High memory capa

city and multi-

GPU support (like NVLink or InfiniBand) are essential for handling large datasets and complex models.

2. **Scalability and Flexibility:** Many cloud providers allow easy scaling to multiple GPUs or nodes, which is crucial for distributed training. Platforms like Runpod and Google Cloud enable users to create multi-node clusters for enhanced performance.
3. **Pricing Models:** Transparent, usage-based pricing is important. Most platforms offer pay-as-you-go models, with options for reserved or spot instances that can reduce costs for flexible workloads. For example, Google Cloud GPUs have competitive pricing starting from \$0.35 per GPU.
4. **User Experience:** A user-friendly interface with quick provisioning and integration with AI tools can significantly reduce setup time. Look for platforms that provide pre-configured environments and easy deployment options.
5. **Security and Compliance:** Ensure that the provider implements robust security measures, including encryption and access controls, especially for enterprise applications. Providers like Runpod operate in certified Tier 3+ data centers, ensuring high reliability and compliance.

### Top Cloud GPU Providers for Machine Learning

- **Runpod:** Offers a scalable AI development platform with a wide selection of GPUs, real-time monitoring, and per-second billing. Ideal for both enterprise and community cloud needs.
- **Google Cloud:** Provides high-performance GPUs with flexible pricing and integration with GCP infrastructure, making it suitable for various AI and machine learning tasks.
- **AWS:** Amazon Web Services offers a range of GPU instances tailored for machine learning, with extensive documentation and support for various frameworks.
- **Hyperstack:** Known for its high-speed networking and cost-saving options, Hyperstack supports a variety of AI workloads and provides a developer-friendly environment.

### GPUs

- Models with a significant number of custom PyTorch/JAX operations that must run at least partially on CPUs
- Models with TensorFlow ops that are not available on Cloud TPU (see the list of [available TensorFlow ops](#))

- Medium-to-large models with larger effective batch sizes

Cloud-

based GPU platforms are essential for modern machine learning projects, offering the necessary computational power and flexibility to train complex models efficiently. When choosing a provider, consider factors such as performance, pricing, scalability, and security to find the best fit for your specific needs.

## Introduction to Cloud TPU

Tensor Processing Units (TPUs) are Google's custom-developed, application-specific integrated circuits (ASICs) used to accelerate machine learning workloads. For more information about TPU hardware, see [TPU architecture](#). Cloud TPU is a web service that makes TPUs available as scalable computing resources on Google Cloud.

TPUs efficiently train your models by using hardware designed for performing large matrix operations often found in machine learning algorithms. TPUs have on-chip high-bandwidth memory (HBM) letting you use larger models and batch sizes. TPUs can be connected in groups called slices that scale up your workloads with little to no code changes.

Code that runs on TPUs must be compiled by the accelerator linear algebra (XLA) compiler. [XLA](#) is a just-in-time compiler that takes the graph emitted by an ML framework application and compiles the linear algebra, loss, and gradient components of the graph into TPU machine code. The rest of the program runs on the TPU host machine. The XLA compiler is part of the TPU VM image that runs on a TPU host machine.

When to use TPUs

Cloud TPUs are optimized for specific workloads. In some situations, you might want to use GPUs or CPUs on Compute Engine instances to run your machine learning workloads.

Tensor Processing Units (TPUs) are Google's custom-developed application-specific integrated circuits (ASICs) used to accelerate machine learning workloads. Cloud TPUs allow you to access TPUs from Compute Engine, [Google Kubernetes Engine](#) and [Vertex AI](#).

## TPUs in Google Cloud

You can use TPUs through Cloud TPU VMs, Google Kubernetes Engine, and Vertex AI. The following table lists resources for each Google Cloud service.

Google service	Cloud	Resources
Cloud TPU		<a href="#">Get started with Cloud TPU VMs</a>
Google Engine	Kubernetes	<a href="#">About TPUs in GKE</a>
Vertex AI		<a href="#">Training on Vertex AI with TPUs</a> <a href="#">Use TPUs for online prediction on Vertex AI</a>

## TPUs

- Models dominated by matrix computations
- Models with no custom PyTorch/JAX operations inside the main training loop
- Models that train for weeks or months
- Large models with large effective batch sizes
- Models with ultra-large embeddings common in advanced ranking and recommendation workloads

Cloud TPUs are *not* suited to the following workloads:

- Linear algebra programs that require frequent branching or contain many element-wise algebra operations
- Workloads that require high-precision arithmetic
- Neural network workloads that contain custom operations in the main training loop

## 5. Experiment tracking and model evaluation in cloud based ML and DL

Cloud-based platforms provide robust tools for experiment tracking and model evaluation in machine learning (ML) and deep learning (DL), enabling reproducibility, comparison, and optimization. Key solutions like MLflow integrate seamlessly across AWS SageMaker, Azure ML, and Google Vertex AI.

Experiment tracking and model evaluation in cloud-based ML/DL involve logging hyperparameters, datasets, and code versions to manage iterations, using tools like MLflow, WandB, or cloud-native services (e.g., Azure ML, SageMaker). These tools enable centralized, automated logging and visualization, facilitating model comparison and reproducibility. Model

evaluation involves analyzing metrics, confusion matrices, and ROC curves to select the best model for production.

### **Key Components of Cloud-Based Experiment Tracking**

- \* **Metadata Logging:** Automated logging of hyperparameters, code versions, environment configurations, and metrics (loss, accuracy) using SDKs.
- \* **Artifact Management:** Storing model weights, plots (confusion matrices, ROC curves), and dataset versions.
- \* **Visualization & Comparison:** Using web UIs to compare hundreds of runs, identify top-performing models, and visualize training progress.
- \* **Collaboration:** Allowing team members to share and review experiment results.

### **Tools and Platforms**

- \* **MLflow:** An open-source, widely adopted framework, often integrated with cloud services for managing the entire ML lifecycle.
- \* **Weights & Biases (W&B):** A popular tool for tracking, visualizing, and organizing experiments.
- \* **Neptune.ai:** A metadata store for logging and comparing ML model experiments.
- \* **DVC (Data Version Control):** Provides Git-like functionality to version data, models, and pipelines, enhancing reproducibility.
- \* **Cloud-Native Solutions:** Amazon AWS SageMaker, Azure Machine Learning, and Google Vertex AI, which offer built-in experiment tracking.

### **Model Evaluation Techniques**

- \* **Metrics Analysis:** Evaluating models using metrics relevant to the task, such as accuracy, precision, recall, F1-score, and AUC.
- \* **Visualizations:** Using confusion matrices, ROC curves, and precision-recall curves to analyze performance.
- \* **Validation Set Performance:** Comparing model predictions against validation or test datasets to ensure generalization.
- \* **Model Registry:** Registering the final, high-performing model in a repository for deployment.

### **Benefits of Cloud-Based Tracking**

- \* **Scalability:** Handling thousands of experiments, this is difficult to manage manually.
- \* **Reproducibility:** Ensuring that experiments can be recreated with the same code and data.

- \* **Cost Efficiency:** Using cloud resources to optimize hyperparameter tuning.
- \* **Automation:** Automatically logging metrics, parameters, and code during training.

## **5.Integration of notebooks(jupyter, colab) with cloud storage**

Jupyter notebooks and Google Colab can integrate seamlessly with major cloud storages like AWS S3, Google Cloud Storage (GCS), and Azure Blob Storage using libraries and authentication methods.

### **AWS S3 Integration**

Use the boto3 library in Jupyter or Colab after configuring AWS credentials via IAM keys.

Create an S3 client to upload, download, or list files: `import boto3; s3 = boto3.client('s3')`.

For Colab, install boto3 with pip and set credentials programmatically.

### **Google Cloud Storage**

Mount GCS buckets directly in Vertex AI JupyterLab via the sidebar or use `google.cloud.storage` client in code.

Authenticate with `gcloud auth login` or `google.colab.auth`, then copy files using `gsutil cp`.

Extensions like `jgscm` enable full Jupyter filesystem support for GCS buckets.

### **Azure Blob Storage**

Connect via `azure-storage-blob` library; define account name, key, and container in Jupyter code.

In Colab, use packages like `mount-azure-blob` to mount containers interactively.

Upload blobs with `BlobServiceClient` after setting connection strings.

### **Google Colab Specifics**

Mount Google Drive natively: `from google.colab import drive; drive.mount('/content/drive')` for easy access.

For cross-cloud, authenticate first (e.g., `google.colab.auth` for GCS) then use client libraries or `gsutil`.

This supports your AWS SageMaker and Vertex AI workflows without local downloads.

Integrating notebooks like Jupyter and Google Colab with cloud storage simplifies data management by providing a persistent, scalable environment for your files.

### **1. Google Colab Integrations**

Colab is designed to live in the Google ecosystem, offering seamless links to storage:

**Google Drive:** You can mount your entire Drive as a local directory using `from google.colab import drive; drive.mount('/content/drive')`. This allows you to read/write files as if they were on your hard drive.

**Google Cloud Storage (GCS):** For larger datasets, use the `google.colab.auth` module to authenticate and then interact with GCS buckets using the Python Client Library or `gsutil` commands.

**GitHub:** Colab can open notebooks directly from a GitHub repository and save them back as commits, making version control straightforward.

## 2. Jupyter & JupyterLab Integrations

Standard Jupyter instances (local or cloud-hosted) use different methods depending on the provider:

**Vertex AI Workbench:** These Google Cloud-hosted JupyterLab instances include a built-in Cloud Storage integration button. You can mount GCS buckets directly to the file browser sidebar to work with files without manual downloading.

**AWS S3:** Use the `s3fs` library to interact with S3 buckets directly within your notebook cells. For example, `pd.read_csv('s3://bucket-name/file.csv')` allows pandas to read data directly from the cloud.

**Azure Blob Storage:** Instances like Azure Machine Learning notebooks automatically integrate with your workspace's default storage account.

## 3. Advanced Mounting & Management

**FUSE (Filesystem in Userspace):** Tools like GCS Fuse allow you to mount cloud buckets as a file system on your VM. This creates a "single source of truth" for multiple data scientists working on the same project.

**LakeFS:** This open-source platform adds Git-like versioning (branching, committing) to your object storage (S3, GCS, Azure Blob), which you can then manage via a Jupyter extension.

## 1. Containers and Docker for AI Applications :

Containers and Docker are essential tools for developing and deploying AI applications in cloud computing. They provide a consistent and scalable environment for AI workloads, ensuring that models and APIs behave identically across different environments. In cloud computing, AI applications require complex software environments, including machine learning libraries, frameworks, and hardware acceleration support. Containers and Docker provide an efficient way to package, deploy, and manage AI applications in the cloud.

**Containers:** A **container** is a lightweight, standalone, and executable software package that includes:

- application code
- libraries and dependencies
- runtime environment
- system tools and configuration

Containers run on a shared operating system kernel but are isolated from each other, making them more efficient than virtual machines.

### Features of Containers

- Lightweight and fast startup
- Platform independent
- High resource efficiency
- Easy scalability and portability

**Docker:** **Docker** is an open-source containerization platform used to create, deploy, and run containers.

Docker uses:

- **Dockerfile** – a script that defines how a container image is built
- **Docker Image** – a packaged application environment
- **Docker Container** – a running instance of an image

### Role of Docker in AI Applications

Docker simplifies AI application deployment by:

- packaging AI models with libraries like TensorFlow and PyTorch
- ensuring consistent environments across development and production
- supporting GPU acceleration using CUDA-enabled images

## Containers in Cloud Computing for AI

In cloud computing, containers enable AI applications to:

- run on different cloud platforms without modification
- scale automatically based on workload
- deploy faster using CI/CD pipelines
- support distributed training and inference

Cloud providers use containers as the foundation for AI services.

### Advantages of Containers and Docker for AI

1. **Reproducibility** – same environment for training and deployment
2. **Scalability** – easy horizontal scaling in the cloud
3. **Portability** – run anywhere (local, cloud, hybrid)
4. **Efficient Resource Utilization** – reduced overhead compared to VMs
5. **Faster Deployment** – quick model updates and rollbacks

### Docker and Kubernetes

- **Docker** is used to build and run containers
- **Kubernetes** is used to manage and orchestrate multiple containers in the cloud

Kubernetes enables load balancing, auto-scaling, and fault tolerance for AI workloads.

### Use Cases

- AI model training
- Machine learning inference services
- Natural language processing applications
- Computer vision systems
- Large-scale data processing

Containers and Docker play a vital role in deploying AI applications in cloud computing. They provide a scalable, portable, and efficient solution for managing complex AI environments, making them essential for modern AI systems. Docker is a powerful tool for containerizing machine learning (ML) projects, ensuring consistency, reproducibility, and simplified deployment.

Below is a step-by-step guide to containerize and run ML models using Docker.

## 1. Set Up Your ML Project

Organize your project directory with the following structure:

my-ml-project/

├── train\_model.py # Script to train the model

├── requirements.txt # Dependencies (e.g., scikit-learn, pandas)

├── data/ # Dataset folder

└── model.pkl # Trained model output

## 2. Create a Dockerfile

The Dockerfile defines the environment and steps to build your container:

```
# Use Python base image
```

```
FROM python:3.9-slim
```

```
# Set working directory
```

```
WORKDIR /app
```

```
# Copy dependencies and install them
```

```
COPY requirements.txt ./
```

```
RUN pip install --no-cache-dir -r requirements.txt
```

```
# Copy project files
```

```
COPY .
```

```
# Run the training script
```

```
CMD ["python", "train_model.py"]
```

## 3. Build the Docker Image

Run the following command in your terminal to build the image:

```
docker build -t my-ml-project
```

## 4. Run the Container

Execute the container to train your model:

```
docker run my-ml-project
```

This will run `train_model.py` inside the container and save outputs like `model.pkl`.

### 5. Persist Data Using Volumes

To save outputs (e.g., trained models) outside the container:

```
docker run -v $(pwd)/output:/app/output my-ml-project
```

Ensure your script saves files to `/app/output`.

### 6. Use Jupyter Notebooks (Optional)

To work interactively with Jupyter:

- Update Dockerfile:

```
RUN pip install jupyterlab
```

```
EXPOSE 8888
```

```
CMD ["jupyter", "lab", "--ip=0.0.0.0", "--port=8888", "--allow-root"]
```

- Build and run:

```
docker build -t ml-jupyter
```

```
docker run -p 8888:8888 ml-jupyter
```

- Access Jupyter at <http://localhost:8888>.

### 7. GPU Support (Optional)

For GPU acceleration:

- Install NVIDIA drivers and NVIDIA Container Toolkit.
- Run with GPU support:

```
docker run --gpus all my-ml-project.
```

## 2. Kubernetes for ai applications in cloud computing

**Kubernetes (K8s)** is an open-source container orchestration platform used to deploy, manage, and scale containerized applications. In cloud computing, Kubernetes plays a crucial role in running **AI and machine learning applications** efficiently at scale.

Kubernetes is a powerful platform for managing AI applications in cloud computing. It enables scalable, reliable, and efficient deployment of AI workloads by orchestrating containerized applications across cloud infrastructure.

## **Why Kubernetes is Needed for AI**

AI applications in the cloud involve:

- large models
- high computational requirements (CPU/GPU)
- variable workloads
- continuous updates of models

Kubernetes automates the management of these complex AI workloads.

## **Key Kubernetes Components Used in AI**

### **1. Pods**

- Smallest deployable unit in Kubernetes
- Contains one or more AI containers
- Runs model training or inference code

### **2. Nodes**

- Physical or virtual machines in the cloud
- Can be CPU or GPU enabled
- Execute AI workloads

### **3. Deployments**

- Manage multiple replicas of AI inference containers
- Ensure high availability

### **4. Services**

- Expose AI models as network services
- Enable load balancing across containers

### **5. ConfigMaps and Secrets**

- Store model configurations and credentials securely

## **Role of Kubernetes in AI Training**

- Schedules training jobs on GPU nodes
- Supports distributed training
- Restarts failed training containers automatically
- Efficient resource allocation

### **Advantages of Kubernetes for AI Applications**

1. **Scalability** – automatic scaling of AI services
2. **High Availability** – self-healing and fault tolerance
3. **Efficient Resource Management** – CPU/GPU scheduling
4. **Portability** – cloud-agnostic deployment
5. **Automation** – simplified deployment and updates

### **Use Cases**

- Large-scale AI model inference
- Distributed deep learning training
- Computer vision and NLP services
- Recommendation systems
- MLOps pipelines

### **3. What Is Serverless Computing: Aws Lambda, Azure Functions.**

Serverless computing is a cloud development model where the provider manages the underlying infrastructure, allowing developers to build and run applications without provisioning or managing servers. AWS Lambda and Azure Functions are the leading Function-as-a-Service (FaaS) platforms within this model offered by Amazon and Microsoft, respectively.

In the serverless model, developers focus solely on writing code (functions) that are executed in response to specific events or requests. **Serverless computing** is a cloud computing model where the cloud provider automatically manages the server infrastructure.

Developers only focus on writing and deploying code, while the cloud provider handles:

- Server provisioning
- Scaling
- Maintenance
- Fault tolerance
- Monitoring

Even though it is called *serverless*, servers still exist, but they are completely managed by the cloud provider.

### Key Characteristics of Serverless Computing

1. **No Server Management** – No need to provision or maintain servers.
2. **Event-Driven Execution** – Code runs in response to events (HTTP requests, file uploads, database updates).
3. **Auto Scaling** – Automatically scales up or down based on demand.
4. **Pay-per-Use Pricing** – Charged only for execution time and resources used.
5. **High Availability** – Built-in fault tolerance.
6. **Short-lived Functions** – Functions run for short duration and terminate.

### Architecture of Serverless Computing

#### Flow of Execution:

Client Request → API Gateway / Trigger → Function Execution → Database / Storage → Response to Client

Example:

- User uploads file → Storage event triggers function → Function processes file → Output stored.

### AWS Lambda

Amazon Web Services provides **AWS Lambda**, a serverless compute service that runs code without managing servers.

**AWS Lambda** allows execution of backend code in response to events.

#### Features of AWS Lambda

1. Fully managed service
2. Supports multiple languages:
  - Python
  - Java
  - Node.js
  - C#
  - Go
3. Automatic scaling
4. Integrated with AWS services
5. Pay per execution (milliseconds billing)

## Working of AWS Lambda

### Step-by-Step Working:

1. Developer uploads function code.
2. Define trigger (e.g., API Gateway, S3, DynamoDB).
3. Event occurs.
4. Lambda executes code.
5. Result returned or stored.

### Example Use Case:

- Image upload in S3 → Lambda resizes image → Stores resized image.

### Components of AWS Lambda

1. **Event Source (Trigger)** – S3, API Gateway, CloudWatch.
2. **Function Code** – Business logic.
3. **Execution Role (IAM Role)** – Permissions to access AWS resources.
4. **Environment Variables** – Configuration settings.

### Advantages of AWS Lambda

- No infrastructure management
- Highly scalable
- Cost efficient
- Easy integration with AWS ecosystem

## Azure Functions

Microsoft Azure provides **Azure Functions**, a serverless compute service that runs event-driven code without managing infrastructure.

### Features of Azure Functions

1. Fully managed serverless platform
2. Supports languages:
  - C#
  - Java
  - Python
  - JavaScript
  - PowerShell
3. Automatic scaling
4. Consumption-based pricing
5. Tight integration with Azure services

## Working of Azure Functions

### Step-by-Step Process:

1. Developer writes function.
2. Define trigger (HTTP, Timer, Blob storage, Queue).
3. Event occurs.
4. Function executes.
5. Output stored or returned.

### Example:

- HTTP request → Azure Function processes data → Saves to Azure SQL Database.

### Types of Triggers in Azure Functions

1. HTTP Trigger
2. Timer Trigger
3. Blob Storage Trigger
4. Queue Trigger
5. Event Hub Trigger

### Advantages of Azure Functions

- Seamless integration with Azure services
- Automatic scaling
- Supports Durable Functions (workflow orchestration)
- Cost-effective

## **4. WHAT IS CI/CD Pipelines for AI Models in Cloud.**

CI/CD pipelines for AI models in the cloud are automated workflows—using tools like AWS CodePipeline or Cloud Build—that streamline the building, training, testing, and deployment of machine learning models. They enable rapid, consistent, and automated updates from code repository to production, allowing teams to quickly retrain, validate, and deploy AI, reducing manual errors and improving accuracy.

### 2. Need for CI/CD in AI Models

AI systems differ from traditional software because:

- Models depend on data quality
- Models require retraining
- Performance may degrade over time (model drift)
- Deployment involves model artifacts

Therefore, automated pipelines ensure:

- Faster experimentation
- Reproducibility
- Version control
- Continuous improvement

### 3. CI/CD Pipeline Stages for AI in Cloud

#### 1. Source Control

- Store code, data schemas, and model configuration in Git.
- Version control of datasets and models.

#### 2. Continuous Integration (CI)

- Automatic build and test when code changes.
- Validate:
  - Data quality
  - Model training scripts
  - Unit tests

#### 3. Model Training

- Training executed on cloud services.
- Model artifacts stored in model registry.

#### 4. Model Validation

- Evaluate accuracy, precision, recall.
- Compare with previous version.
- Approve or reject model.

#### 5. Continuous Deployment (CD)

- Deploy model to:
  - API endpoints
  - Batch processing systems
  - Edge devices

#### 6. Monitoring

- Monitor:
  - Accuracy
  - Latency
  - Data drift
- Trigger retraining if required.

## 4. Architecture of AI CI/CD Pipeline

### Flow Diagram:

Code Commit → CI Build → Data Validation → Model Training → Model Evaluation → Model Registry → Deployment → Monitoring → Retraining

## CI/CD Tools in Major Cloud Platforms

### 5. CI/CD in AWS for AI

Amazon Web Services provides AI CI/CD through:

- **Code Pipeline** – Automates workflow
- **CodeBuild** – Build and test
- **SageMaker Pipelines** – ML workflow automation
- **ECR** – Container registry

Typical Flow:

GitHub → CodePipeline → CodeBuild → SageMaker Training → Model Registry → Deployment to Endpoint

### 6. CI/CD in Azure for AI

Microsoft Azure supports:

- **Azure DevOps Pipelines**
- **Azure Machine Learning Pipelines**
- **Azure Container Registry**

Flow:

Code Commit → Azure DevOps → Azure ML Training → Model Validation → Deploy as REST Endpoint

### 7. CI/CD in GCP for AI

Google Cloud Platform provides:

- Cloud Build
- Vertex AI Pipelines
- Artifact Registry

Flow:

Source Repo → Cloud Build → Vertex AI Training → Model Registry → Deployment

## 8. Advantages of CI/CD for AI Models

1. Faster deployment cycles
2. Improved collaboration
3. Automated testing
4. Reduced manual errors
5. Continuous model improvement
6. Better monitoring and governance

## 9. Challenges in AI CI/CD

1. Data versioning complexity
2. Large training time
3. Model reproducibility issues
4. Managing model drift
5. Infrastructure cost

## 5. Scaling AI Applications using Load Balancers and Auto-Scaling.

Scaling AI applications requires combining Auto-Scaling for dynamic capacity adjustment and Load Balancers for traffic distribution to ensure high availability, performance, and cost efficiency. Auto-scaling adds/removes instances based on CPU/GPU utilization or request rates. Load balancers distribute incoming requests across these instances, directing traffic away from unhealthy units.

AI applications such as chatbots, recommendation systems, fraud detection systems, and image recognition services often receive unpredictable and heavy traffic. To maintain performance, availability, and low latency, these applications must scale dynamically.

Scaling is achieved using:

- **Load Balancers**
- **Auto-Scaling Mechanisms**

These ensure efficient resource utilization and uninterrupted service.

## 2. Need for Scaling in AI Applications

AI systems require scaling because:

1. High computational requirements (GPU/CPU intensive).
2. Variable user traffic.
3. Real-time inference requirements.
4. Large data processing workloads.
5. High availability expectations.

Without proper scaling, systems may experience latency, failures, or downtime.

### 3. Load Balancing in AI Applications

#### Definition

A **Load Balancer** distributes incoming network traffic across multiple servers or instances to ensure no single server is overloaded.

#### Working of Load Balancer

Client Requests → Load Balancer → Multiple AI Model Instances → Response

The load balancer:

- Monitors server health
- Routes traffic to healthy instances
- Ensures even distribution

#### Types of Load Balancers

1. **Application Load Balancer (Layer 7)** – HTTP/HTTPS based routing.
2. **Network Load Balancer (Layer 4)** – TCP/UDP traffic routing.
3. **Global Load Balancer** – Multi-region distribution.

#### Example in Cloud

- Amazon Web Services provides **Elastic Load Balancer (ELB)**.
- Microsoft Azure provides **Azure Load Balancer** and Application Gateway.
- Google Cloud Platform provides **Cloud Load Balancing**.

### 4. Auto-Scaling in AI Applications

#### Definition

**Auto-Scaling** automatically increases or decreases the number of compute instances based on workload demand.

#### Working of Auto-Scaling

1. Monitor metrics (CPU usage, memory, request count).
2. If usage exceeds threshold → Add instances (Scale Out).
3. If usage decreases → Remove instances (Scale In).

## Types of Scaling

1. **Vertical Scaling (Scale Up)** – Increase CPU/RAM of a machine.
2. **Horizontal Scaling (Scale Out)** – Add more machines/instances.

## Cloud Auto-Scaling Services

- AWS Auto Scaling
- Azure Virtual Machine Scale Sets
- GCP Managed Instance Groups

## 5. Combined Architecture for AI Scaling

### Diagram Explanation (Text Format)

User Requests

↓

Load Balancer

↓

Auto-Scaling Group

↓

Multiple AI Model Instances

↓

Database / Storage

If traffic increases → Auto-scaling adds more instances.  
Load balancer distributes traffic among them.

## 6. Advantages of Load Balancing + Auto-Scaling

1. High availability
2. Fault tolerance
3. Reduced latency
4. Cost optimization
5. Better performance during peak loads
6. Efficient resource utilization

## 7. Challenges

1. Cold start in AI models
2. GPU resource cost
3. Model synchronization across instances
4. Data consistency issues.

## Example: AI-Based E-Commerce Recommendation System

Consider an e-commerce company like Amazon that uses an AI recommendation system to suggest products to users.

The AI model analyzes:

- User browsing history
- Purchase behavior
- Search queries
- Product preferences

During events like **Big Billion Sale** or festival sales, millions of users access the platform simultaneously. The AI system must handle huge traffic without delay.

## 6. Monitoring and Logging in Cloud for AI Workflows.

Monitoring and Logging are essential components of AI workflows deployed in cloud environments.

- **Monitoring** tracks the performance, health, and behavior of AI models and infrastructure.
- **Logging** records detailed event information for debugging, auditing, and analysis.

Monitoring and logging are essential for maintaining the reliability, performance, and compliance of AI workflows in the cloud. As AI systems scale, traditional infrastructure monitoring is often supplemented by **AI observability**, which tracks model-specific behaviors like data drift and inference accuracy alongside standard system metrics.

### Key Components of AI Cloud Monitoring

- **Infrastructure Metrics:** Continuous tracking of CPU, GPU, memory, and network utilization to ensure the environment can handle heavy AI computation.
- **Model Performance:** Monitoring metrics such as inference latency, prediction accuracy, precision, recall, and F1 score to detect slips in model quality.
- **Data & Model Drift:** Identifying when real-world production data deviates from the training set (**data drift**) or when the model's predictive power degrades over time (**model drift**).
- **Cost Tracking:** Monitoring token consumption and cloud resource spend to prevent budget overruns in expensive AI operations.

### Logging Strategies for AI Workflows

Cloud-native logging provides the raw data needed to trace anomalies and maintain audit trails for regulated industries.

- **Execution Logs:** Automatically generated records of workflow starts, ends, and transitions between steps.
- **Input/Output Logs:** Capturing prompts and responses (redacting sensitive data) to analyze agent reasoning and identify hallucinations.

- **Distributed Tracing:** Using traces and spans (e.g., via **OpenTelemetry**) to map requests across multi-stage AI pipelines, such as data retrieval, model inference, and post-processing.

## Cloud Tools for Monitoring AI

### AWS

Amazon Web Services provides:

- CloudWatch (monitoring)
- CloudTrail (logging & auditing)
- SageMaker Model Monitor

### Azure

Microsoft Azure provides:

- Azure Monitor
- Application Insights
- Azure ML Monitoring

### GCP

Google Cloud Platform provides:

- Cloud Monitoring
- Cloud Logging
- Vertex AI Model Monitoring

## Benefits of Monitoring and Logging in AI

1. Early detection of model drift
2. Improved reliability
3. Faster troubleshooting
4. Regulatory compliance
5. Improved decision-making
6. Continuous model improvement

## UNIT IV

### 1. Security and Privacy in Cloud-based AI Ethical

Securing Cloud-based AI requires a delicate balance between leveraging vast computational power and enforcing strict **data privacy** and **ethical standards**. Key risks include algorithmic bias, model theft, and unintentional data exposure, all of which demand robust **Privacy by Design** and continuous **AI Security Posture Management (AI-SPM)**.

#### 1. Ethical Challenges in Cloud AI

- **Algorithmic Bias & Fairness:** AI models trained on unrepresentative cloud datasets can reinforce discrimination. It is crucial to employ fairness-aware modeling to ensure equitable outcomes.
- **Transparency & Explainability:** Advanced cloud models often operate as "black boxes". Using **Explainable AI (XAI)** helps stakeholders interpret AI decisions, fostering trust and accountability
- **Consent & Governance:** Repurposing personal data or web-scraped information for AI training presents severe ethical and copyright hurdles. Organizations must establish clear data governance and consent policies.

#### 2. Privacy & Data Protection

- **Data Minimization:** Only collect and process personal data that is strictly necessary for training or inference tasks.
- **Privacy-Preserving Technologies:** Utilize cryptographic and analytical techniques such as **Differential Privacy** to allow models to learn from data patterns without accessing identifiable, raw information.
- **Anonymization & De-identification:** Strip Personally Identifiable Information (PII) from datasets before feeding them into cloud environments.

#### 3. Security Vulnerabilities & Defense

- **Cloud & Model Leaks:** Generative AI solutions can accidentally leak confidential company data if not properly sandboxed in the cloud.
- **Data Poisoning & Adversarial Attacks:** Attackers may manipulate or alter training datasets to corrupt AI outputs.
- **Model Theft:** Sensitive AI models stored in multi-tenant cloud architectures can be stolen or reverse-engineered.

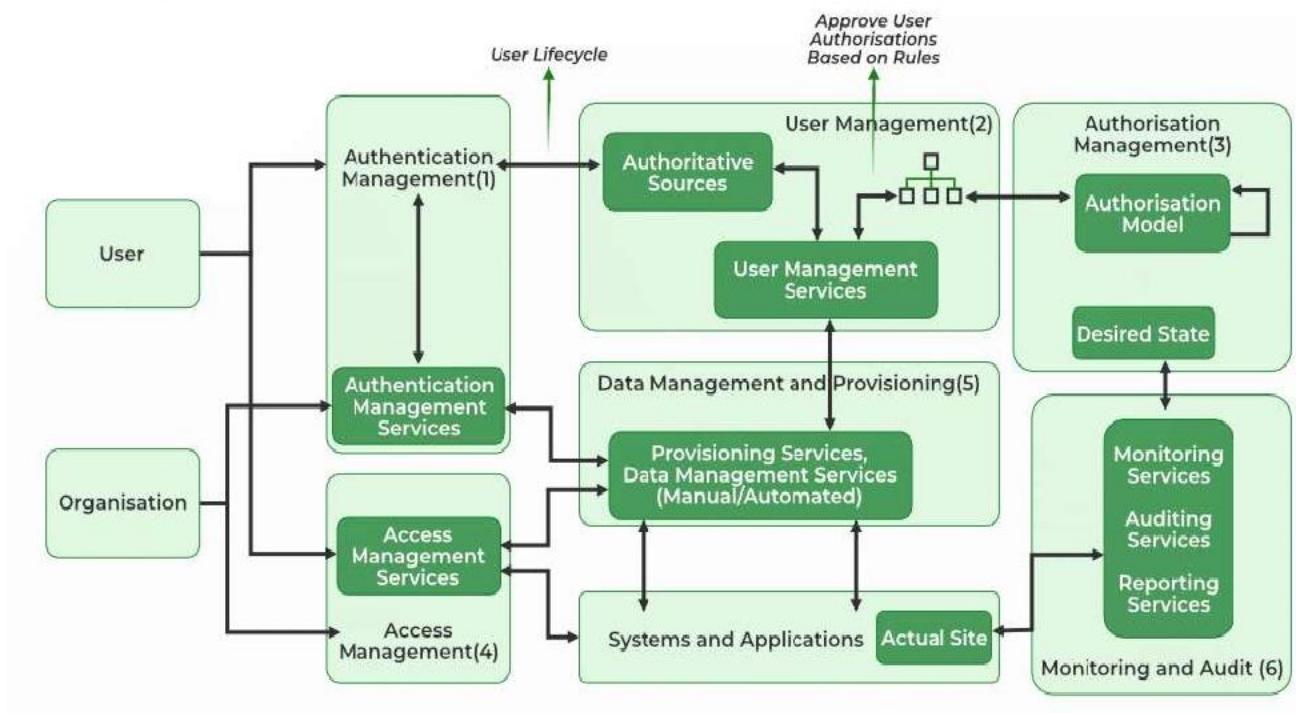
- **Zero-Trust Framework:** Implement zero-trust architectures that never assume user trust to protect against unauthorized access and ensure compliance with regulations like the **EU AI Act** and **GDPR**.

## 2. Identity and Access Management (IAM) in Cloud

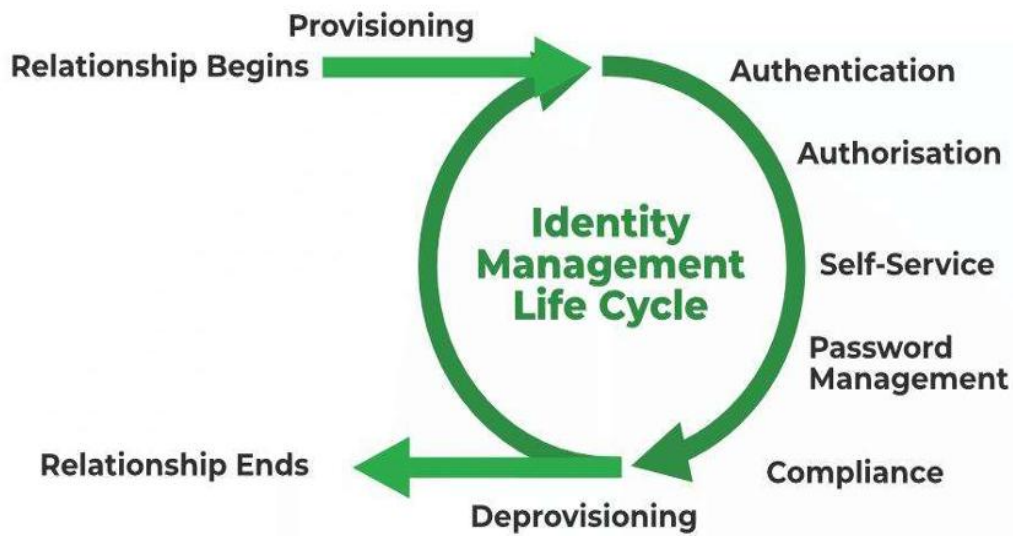
IAM architecture is the framework that helps securely manage **user identities, authentication** and **authorization** across an organization's IT environment.

**The Architecture of IAM in cloud computing:**

- **User Management:** It consists of activities for the control and management over the identity life cycles.
- **Authentication Management:** It consists of activities for effectively controlling and managing the processes for determining which user is trying to access the services and whether those services are relevant to him or not.
- **Authorization Management:** It consists of activities for effectively controlling and managing the processes for determining which services are allowed to access according to the policies made by the administrator of the organization.
- **Access Management:** It is used in response to a request made by the user wanting to access the resources with the organization.

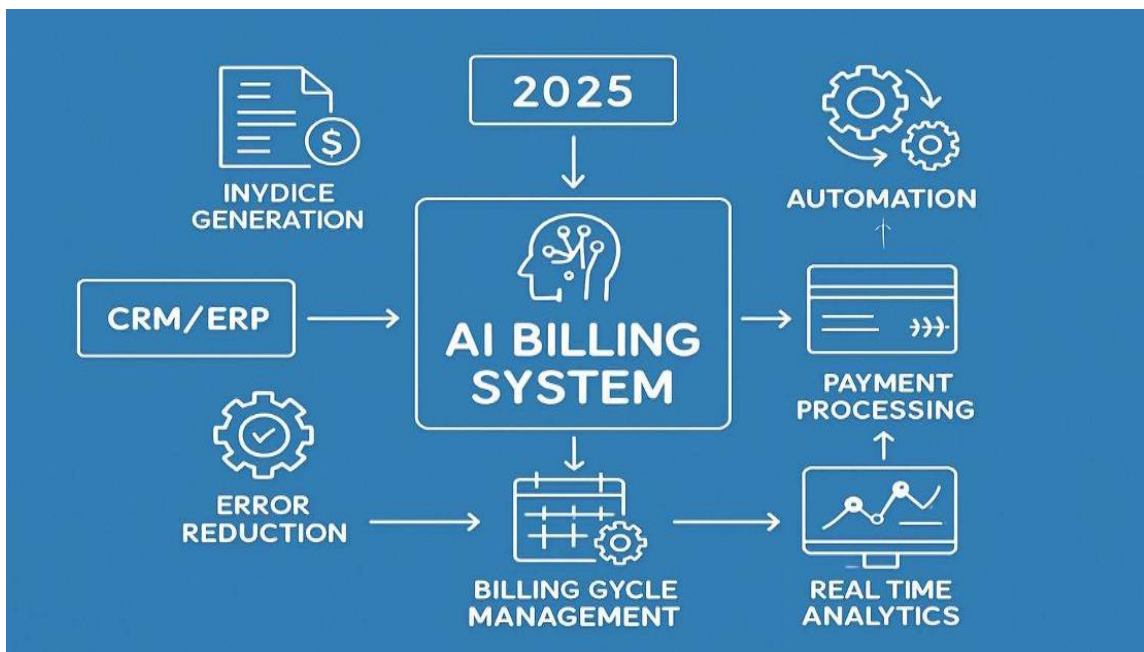


- **Data Management and Provisioning:** The authorization of data and identity are carried towards the IT resource through automated or manual processes.
- **Monitoring and Auditing:** Based on the defined policies the monitoring, auditing, and reporting are done by the users regarding their access to resources within the organization.
- **Operational Activities of IAM:** In this process, we onboard the new users on the organization's system and application and provide them with necessary access to the services and data. Deprovisioning works completely opposite in that we delete or deactivate the identity of the user and de-relinquish all the privileges of the user.
- **Credential and Attribute Management:** Credentials are bound to an individual user and are verified during the authentication process. These processes generally include allotment of username, static or dynamic password, handling the password expiration, encryption management, and access policies of the user.
- **Entitlement Management:** These are also known as authorization policies in which we address the provisioning and de-provisioning of the privileges provided to the user for accessing the databases, applications, and systems. We provide only the required privileges to the users according to their roles. It can also be used for security purposes.
- **Identity Federation Management:** In this process, we manage the relationships beyond the internal networks of the organization that is among the different organizations. The federations are the associate of the organization that came together for exchanging information about the user's resources to enable collaboration and transactions.
- **Centralization of Authentication and Authorization:** It needs to be developed in order to build custom authentication and authorization features into their application, it also promotes the loose coupling architecture.



### 3. Cost management and billing for Ai services

AI cost management involves tracking, optimizing, and forecasting usage-based expenses like tokens, compute seconds, and API calls. Because AI usage is highly unpredictable, effective management relies on setting strict hard caps, utilizing continuous anomaly detection, and optimizing application architectures for better cost-to-business outcomes.



## Key AI Pricing Models

- **Usage-Based (Inference):** Billed per input/output token (text), audio duration, or image resolution/generation.
- **Dedicated Compute (Training/Hosting):** Billed via continuous hourly rates for running dedicated virtual machines or fine-tuning environments.
- **Credit-Based:** Pre-allocated pools of funds (e.g., [Microsoft Copilot Credits](#)) that are depleted as end-users consume services.

## Best Practices for Billing Control

- **Set Hard Limits & Budgets:** Avoid surprise 10x bills by setting proactive limits on your usage-based monthly spending.
- **Track Anomalies:** Treat unexpected cost spikes as high-priority incidents and address them with engineering before the month ends.
- **Normalize Business Metrics:** Instead of just measuring raw tokens, calculate your effective cost per business outcome (e.g., cost per 1,000 document summaries).
- **Isolate Environments:** Keep production, staging, and continuous-running development/testing environments billed separately to pinpoint waste.

## Essential AI Management Tools

- **Azure Cost Management / AWS / GCP:** Major cloud providers offer built-in analytical assistants where you can type natural language prompts (e.g., *"Why is my AI bill higher this month?"*) to get instant optimization suggestions.
- **Dedicated AI Platforms:** Providers like [Render](#) offer fixed monthly pricing and autoscaling limits for AI applications.
- **Third-Party Platforms:** Tools such as Cast AI, Kubecost, and CloudCheckr help track and automate resource rightsizing across multi-cloud setups.

## 4. Ethical Issues and Fairness in Cloud AI

Ethical issues in Cloud AI largely revolve around algorithmic bias, data privacy, transparency, and accountability. Because cloud-deployed models rely on large datasets, they can inadvertently perpetuate societal biases or function as "black boxes," requiring strict governance and continuous mitigation.

### Major Ethical & Fairness Challenges



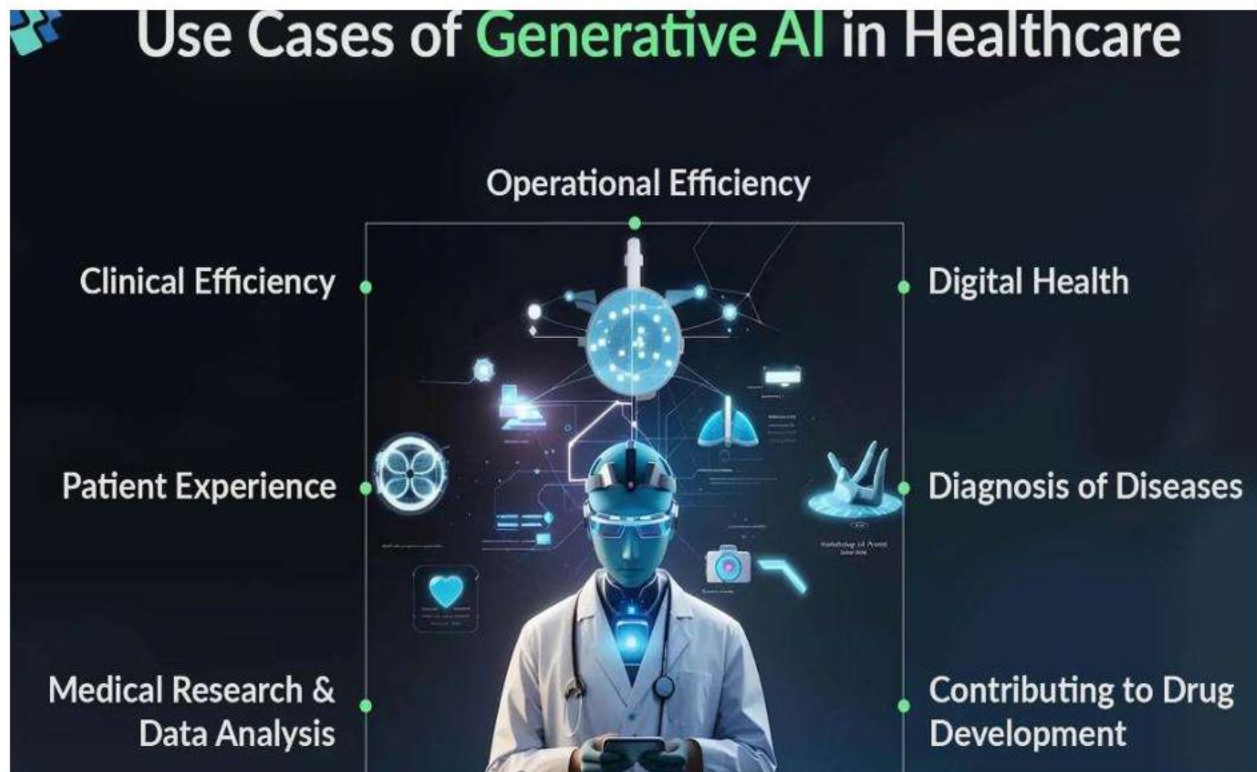
- **Algorithmic Bias:** Machine learning models learn from historical datasets. If the data is skewed or contains historical prejudices, the AI will likely produce discriminatory or unjust outcomes (e.g., in recruitment, lending, or law enforcement).
- **Lack of Explainability:** Many complex AI models deployed on cloud platforms function as "black boxes". It can be difficult to interpret *why* a specific decision was made, making it hard to identify and correct errors.
- **Data Privacy & Consent:** Cloud environments require vast amounts of data. Safeguarding sensitive personal information and ensuring users give informed consent are foundational ethical requirements.
- **Accountability:** When an automated cloud system makes an error—especially one that impacts individuals in healthcare or finance—determining who is responsible (developer, cloud provider, or user) is complex.

### Mitigation Strategies

- **Bias Detection and Mitigation:** Organizations can use fairness toolkits (like IBM's AI Fairness 360) to detect and correct algorithmic biases during the model training process.
- **Explainable AI (XAI):** Developers can use specific frameworks to interpret model predictions, turning "black box" decisions into understandable, transparent outcomes.
- **Representative Data Collection:** Ensuring that training datasets are diverse, inclusive, and rigorously audited before deployment helps prevent skewed outcomes across different demographic groups.
- **Comprehensive Governance:** Technical solutions alone are insufficient. Implementing broad governance frameworks, continuous system monitoring, and periodic audits are vital for maintaining system fairness over time.

## 5. Case Study: AI in Healthcare Cloud Solutions

AI –powered cloud solutions are transforming the health care sector by automating workflows, reducing costs , and improving patient outcomes. Organizations leverage scalable cloud platforms to securely process massive medical datasets and deploy real-time clinical tools.



Real-world case studies demonstrate exactly how these cloud and AI integrations operate in practice:

### 1. Clinical Documentation and Workflow Automation

- **The Solution:** **AWS HealthScribe** and **Amazon Transcribe** are leveraged by healthcare providers to automate medical transcripts and notes.
- **The Impact:** These AI solutions recognize medical speech and classify dialogues to generate preliminary clinical transcripts in real-time, drastically reducing the administrative burden on doctors.
- **Read more:** Explore this **AWS & AI in Healthcare Case Study** by JetBase to see how millions of medical device readings are processed for actionable insights.

## 2. Operations Intelligence and Mobile AI Assistants

- **The Solution:** Providers utilize the **Microsoft Azure** AI platform to build natural language processing (NLP) operational assistants and touch-optimized mobile apps.
- **The Impact:** Healthcare managers and nurses can interactively query enterprise data to handle ICU staffing alerts, predict operational needs, and process complex analytical queries in less than two minutes directly from mobile devices.
- **Read more:** Review this Enterprise AI Solution Case Study by Diaspark to learn about optimizing healthcare personnel workflows.

## 3. Cost Reduction via Cloud Transformation

- **The Solution:** Hospitals migrate vast historical mainframe databases and electronic health records (EHR) to scalable, HIPAA-compliant cloud architectures (e.g., Azure SQL Servers and Tableau).
- **The Impact:** This eliminates the need for expensive on-premises hardware, saving organizations up to \$1 million annually in mainframe costs while maintaining regulatory compliance.
- **Read more:** See the Health Care Cloud Migration Case Study by Deloitte for a deeper dive into mainframe modernization.

## 4. Patient Engagement and Analytics

- **The Solution:** Consultancies like **Cognizant** deploy AI assistants like RESOLV using cloud-based enterprise data and CRM systems.
- **The Impact:** These assistants provide narrative interpretations of medical and business charts, and even help alert doctors to at-risk patient behaviors during visits to promote preventive care.
- **Read more:** Check out the Cognizant AI Healthcare Case Study to see how smart assistants simplify complex analytics.

## 6. Case Study: Real-Time Analytics in Financial Cloud Services

Real-time cloud analytics in finance enables institutions to process transactions and market data instantly. By utilizing high-speed ingestion and cloud-based databases, firms can reduce fraud by 30%, eliminate manual spreadsheet compiling, and increase forecasting accuracy.

## Key Use Cases



- **Fraud Detection & AML:** Streaming analytics process millions of transactions per second. Using tools like [Striim](#) or [Apache Kafka](#), machine learning models score transactions instantly. This allows banks to flag and block fraudulent behavior or money laundering in milliseconds, preventing monetary losses. **Live Financial Reporting:** Cloud solutions allow finance departments to transition from monthly, static Excel reports to live, interactive dashboards. This reduces administrative reporting time to zero and ensures executive decisions are made with up-to-date data. **Algorithmic & High-Frequency Trading:** Real-time databases enable firms to instantly react to market fluctuations, stream live feeds, and execute orders before market conditions change

- **Core Cloud Architecture**

Building real-time financial services involves a multi-layer pipeline that operates in milliseconds:

1. **Ingestion Layer:** High-velocity transaction streams are captured using streaming frameworks like Apache Kafka or AWS Kinesis.
2. **Processing Layer:** Frameworks like Apache Spark or Storm process live data feeds to filter out noise, deduplicate events, and perform calculations.
3. **Analytics & Storage Layer:** Cloud-native platforms like [CrateDB](#) or Data Warehouses like Snowflake aggregate live data for ML model inferencing.
4. **Action & Visualization:** Machine Learning outputs trigger automated rule engines (e.g., block card) or update [interactive dashboards](#) for risk management.

**Benefits & ROI**

- **Cost Savings & Compliance:** Real-time visibility ensures regulatory compliance and avoids operational delays caused by processing lag.
- **Enhanced Decision-Making:** Access to live financial data allows for rapid reaction to market conditions and improves forecasting accuracy by up to 30%.