

## UNIT IV

### Cloud Computing

**Networking for Cloud Computing:** Introduction, Overview of Data Center Environment, Networking Issues in Data Centers, Transport Layer Issues in DCNs, TCP Enhancements for DCNs, Cloud Service Providers: Google, Amazon Web Services, Microsoft, IBM, SAP Labs, Salesforce, Rackspace, VMware.

#### 1. Introduction to Data Centers

##### What is a Data Center?

A **data center** is a large facility that stores and manages huge amounts of computer systems, servers, and networking equipment.

It acts like the “**brain**” of the Internet, where all online services are processed and stored.

Whenever we:

- Search on **Google**,
- Watch videos on **YouTube**,
- Use **Facebook**, or
- Shop on **Amazon**,

— our requests go to a **data center** somewhere in the world.

##### Types of Data Centers

Data centers can be divided mainly into **two types**:

###### 1. Service-based Data Centers

These provide **online services** directly to users.

**Examples:**

- Google Search
- Facebook
- Yahoo Mail

These companies focus on providing **fast and reliable online experiences** to users.

###### 2. Resource-based Data Centers

These provide **computing resources** (like virtual machines, storage, or processing power) to others.

**Examples:**

- Amazon Web Services (**AWS EC2**)
- Microsoft **Azure**

These allow other people or companies to rent computing power from them instead of buying their own hardware.

## Data Centers and Cloud Computing

Data centers are the **foundation of cloud computing**.

Cloud computing means using the Internet to access servers, storage, and software instead of having everything on your own computer.

So, when we say:

“My data is in the cloud,”

we actually mean

“My data is stored in a data center somewhere.”

Because of cloud computing, many businesses no longer need to buy and maintain their own expensive servers — they can **rent space and services** from big providers like Google Cloud or AWS.

## Importance of Efficient Data Centers

A **good data center** helps to:

- Process data faster
- Keep services running 24/7 without downtime
- Reduce energy and maintenance costs
- Secure user information

But managing such large systems is **complex** and **costly** — that's why only big companies can afford to build and maintain them.

## Tier Classification of Data Centers

To measure how reliable a data center is, they are divided into **four tiers (Tier I to Tier IV)**:

Tier	Description	Example
<b>Tier I</b>	Basic setup with limited backup. Suitable for small businesses.	Small web hosting company
<b>Tier II</b>	Has some backup for power and cooling.	Small IT office
<b>Tier III</b>	Fully redundant systems, 99.982% uptime.	Large corporate data centers
<b>Tier IV</b>	Highest reliability, 99.995% uptime, full redundancy, and fault tolerance.	Banks, hospitals, or large cloud providers like Google and Amazon

## 💰 Cost Distribution in Data Centers

Running a data center is expensive. The major costs usually go into:

- **Servers and hardware**
- **Cooling systems (air conditioning)**
- **Power consumption**
- **Networking equipment**
- **Maintenance and staff**

## 💡 Example:

In a typical cloud data center, around **45% of the cost** may go toward **servers**, **25%** for **power and cooling**, and the rest for **network and operations**.

## 🚀 Modern Cloud Data Centers

Companies like **Google, Microsoft, Facebook, and Amazon** run massive **Internet Data Centers (IDCs)**.

These centers host millions of servers that handle:

- Web search
- File storage (like Google Drive)
- Social networking
- E-commerce (like Amazon)
- Big data processing (like analytics and AI)

## ⚡ Challenges Faced by Data Centers

Even though data centers are powerful, they face some key challenges:

### 1. High Power Consumption 🌩

→ Cooling and running thousands of servers consume a lot of electricity.

### 2. Network Bottlenecks 🌐

→ Fast and efficient data transfer inside the data center is essential.

### 3. Cost Optimization 💰

→ Balancing performance with cost is always difficult.

### 4. Security and Privacy 🔒

→ Protecting user data from cyber threats is a top priority.

## 2. Overview of Data Center Environment

### 1. From Server Rooms to Data Centers

Before modern data centers existed, companies used to have **server rooms** inside their offices. These rooms contained:

- Servers 
- Network switches and routers 
- Cables and backup power supplies 

They helped connect computers in the office through a **Local Area Network (LAN)**.

#### Example:

A bank or college might have one small room with a few servers that stored customer or student data. That was their “server room.”

### 2. Main and Standby Server Rooms

To prevent data loss or downtime, many organizations had:

- **One main server room** (for daily operations)
- **One standby server room** (backup for emergencies)

If something went wrong in the main room — like a **power failure, fire, or hardware issue** — the **standby room** would take over, so the services could continue running without interruption.

#### Example:

If a company’s main server room in Building A went down, its standby server room in Building B could immediately start running the same services.

This system improved **fault tolerance**, meaning the company could still operate even if part of its system failed.

### 3. Rise of Client–Server Computing and the Internet

When **client–server computing** and the **Internet** became popular in the 1990s, the demand for servers grew rapidly.

Companies started needing **more space, more servers, and better cooling and security**.

That’s when the idea of **data centers** was born.

They became **centralized facilities** built specifically to:

- Store servers
- Manage data
- Ensure reliability and uptime

#### Example:

During the **dot-com boom (late 1990s – early 2000s)**, many Internet companies like Yahoo, Google, and Amazon started using large data centers to host their websites and services.

## 🔒 4. Mission-Critical Infrastructure

The term “**mission critical**” means the systems are **essential** for business survival.

A data center is **mission critical** because if it stops working, the company’s entire online service can fail.

That’s why data centers are often described as **fortresses** — they are built for:

- **Maximum reliability**
- **Constant uptime (24/7 operation)**
- **Security and redundancy**

### 💡 Example:

If Amazon’s data center goes offline even for 10 minutes, millions of orders, payments, and deliveries can be disrupted — leading to huge financial losses.

## ⚙️ 5. Evolution and Modern Data Centers

In the last few decades, data centers have **evolved** dramatically:

- **Earlier:** Focus was mainly on reliability — keeping servers running.
- **Now:** Focus includes **efficiency, automation, sustainability, and scalability**.

In the last **5–10 years**, innovation has accelerated due to:

- Cloud computing 
- Virtualization 
- Artificial Intelligence (AI) 
- Green energy 

Modern data centers can automatically balance workloads, reduce energy use, and recover from failures faster than ever.

---

## 2.1 Architecture of Classical Data Centers

### 🌐 What is a Data Center Architecture?

A **data center architecture** is the **blueprint or structure** that shows how servers, storage devices, and networking components are arranged and connected.

It determines how data flows, how fast applications run, and how reliable the system is.

In simple terms:

A data center architecture is like the **layout of a city**, where different buildings (servers) perform different functions but work together through roads (networks).

## Purpose of Data Center Architecture

A well-designed data center ensures:

- **High performance** – fast data access and processing
- **Resiliency** – continues working even if some parts fail
- **Scalability** – can grow easily when demand increases

Proper planning is essential because data centers support all the major operations of a company — from emails and databases to customer services and online transactions.

## Main Types of Classical Data Center Architectures

There are **two main architectural models** commonly used in classical (traditional) data centers:

1. **Multitier Model**
2. **Server Cluster Model**

### **1 Multitier Model**

#### Overview

The **multitier model** is the **most common** data center design.

It divides applications into **multiple layers (tiers)** — each performing a specific task.

This approach is often used for:

- **E-commerce websites** 
- **Enterprise Resource Planning (ERP)** systems
- **Customer Relationship Management (CRM)** systems

It supports technologies such as:

- **Microsoft .NET Framework**
- **Java 2 Enterprise Edition (J2EE)**
- Solutions from **Oracle** and **Siebel**

#### Layers (Tiers) of the Multitier Model

A **three-tier architecture** is most common:

Tier Name	Function / Role	Example
<b>1 Web Server Tier</b>	Handles requests from users through browsers. Apache, Nginx, IIS	
<b>2 Application Server Tier</b>	Processes business logic and rules.	Tomcat, WebLogic
<b>3 Database Server Tier</b>	Stores and manages data.	MySQL, Oracle DB

## How It Works (Example)

Let's say you're shopping on **Amazon**:

1. You open the website → request goes to the **Web Server**.
2. You search for a product → **Application Server** processes your request.
3. The app retrieves product details from the **Database Server** and sends them back to you.

Thus, all three tiers work together to deliver one smooth experience.

## Advantages of Multitier Architecture

-  **Resiliency (Fault Tolerance):**  
If one server fails, others can continue working.
-  **Security:**  
Even if the web server is hacked, attackers cannot directly reach the database.
-  **Scalability:**  
New servers can be added easily at any tier.
-  **Performance:**  
Load can be distributed across multiple servers (called **server farms**).

## Example:

A university's online portal might have:

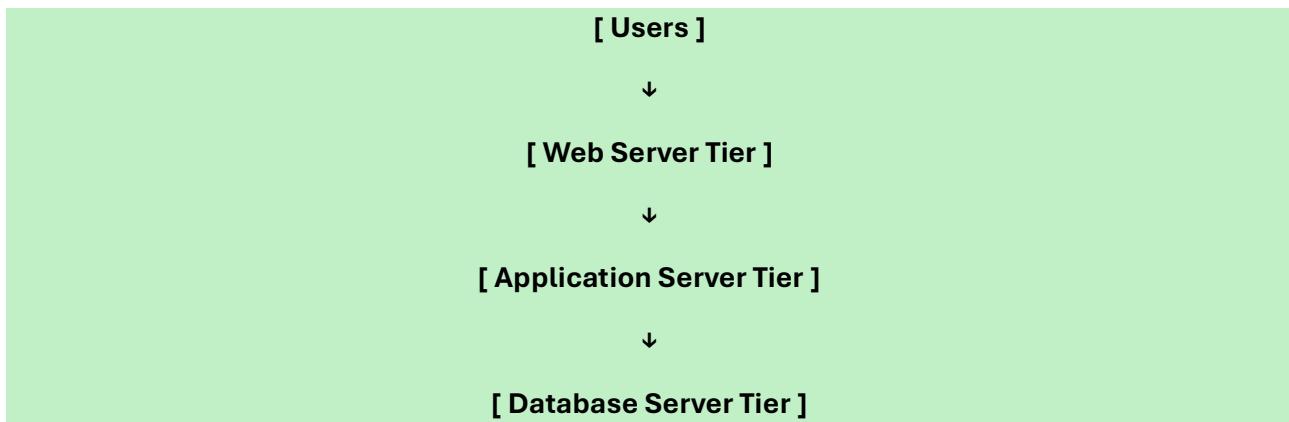
- Web server for students to log in
- Application server for attendance and marks processing
- Database server for storing student records

## Server Farms

In large organizations, each tier may consist of **multiple servers** (called a **server farm**) instead of one single server.

This improves:

- Performance (load sharing)
- Availability (backup servers ready if one fails)



## 2 Server Cluster Model

### Overview

The **server cluster model** originated from **scientific and university research environments** where high computing power was needed for complex calculations.

Over time, it became popular in:

- **Financial institutions**
- **Manufacturing companies**
- **Entertainment (like movie rendering and gaming)**

### How It Works

A **cluster** means a group of servers working together as a single system.

If one server fails, another automatically takes over — making it **highly reliable** and **powerful**.

Used in:

- **High-Performance Computing (HPC)**
- **Parallel Computing**
- **High-Throughput Computing (HTC)**
- **Grid / Utility Computing**

### Example:

NASA or ISRO might use a **cluster of 100 servers** to process space data simultaneously.

### Key Differences: Multitier vs. Server Cluster

Feature	Multitier Model	Server Cluster Model
Purpose	Supports web and business applications	Supports scientific or high-performance computing
Structure	Divided into web, app, and database tiers	Multiple servers work as one powerful system
Examples	E-commerce sites, ERP, CRM	Research labs, weather forecasting, animation rendering
Focus	Data flow and application security	Processing speed and reliability

## 2.2 Cloud-Enabled Data Centers (CEDCs)

A **Cloud-Enabled Data Center (CEDC)** is an **advanced form of a data center** that uses **cloud computing technologies** to become more **flexible, intelligent, and automated**.

In simple terms:

A CEDC = A **modern data center** that uses **cloud principles** (like virtualization, automation, and scalability) to work smarter and faster.

### From Virtualization to CEDC

Before understanding CEDCs, let's recall what **virtualization** means:

Virtualization allows one physical server to act like **many virtual servers**.

Each virtual machine (VM) can run its own applications independently.

But virtualization alone is not enough.

A **CEDC** takes virtualization **to the next level** by adding intelligence, automation, and integration.

### Key Features of CEDCs

#### 1. Virtualization Extended

- CEDCs make better use of virtualization by distributing workloads dynamically.
- Each virtual machine gets the **exact amount of resources (CPU, memory, storage)** it needs.

#### Example:

During online sale days (like Amazon Great Indian Festival), the CEDC automatically increases computing resources for order processing — and reduces them later to save cost.

#### 2. Integration with Cloud Layers

- CEDCs integrate **IaaS (Infrastructure as a Service)** and **PaaS (Platform as a Service)**.
- This means they not only provide servers and storage but also software platforms to run applications efficiently.

#### Example:

A developer can deploy an app using a CEDC without worrying about the hardware or network — the platform handles it automatically.

#### 3. Automation and Orchestration

- **Automation:** Tasks like deploying servers or balancing load happen automatically.
- **Orchestration:** Coordinates these automated tasks across different systems.

#### Example:

If one data center is overloaded, orchestration tools automatically shift some workloads to another data center.

#### 4. Dynamic Workload Management

- CEDCs monitor workloads continuously.
- They allocate resources dynamically — giving more computing power to busy applications and less to idle ones.

##### 💡 Example:

A video streaming site like YouTube automatically gets more bandwidth and servers during peak hours (evenings) and scales down during off-hours.

#### 5. Support for Heterogeneous Data Centers

- “Heterogeneous” means **different types of hardware or platforms**.
- A CEDC can connect and manage **multiple data centers** (from different locations or companies) as one unified system.

### 2.3 Physical Organization of Data Centers

#### 🧠 What Does “Physical Organization” Mean?

The **physical organization** of a data center refers to **how the equipment is arranged inside the building** — including servers, racks, cables, cooling systems, and network connections.

In simple terms:

It's the **layout or physical design** of how everything is placed inside a data center.

#### 🏠 1 Size and Structure of Data Centers

Data centers can vary in size depending on the organization's needs:

- A **small company** may have a data center that fits in a **single room**.
- A **medium-sized company** might use **a few floors of a building**.
- **Big tech companies** like Google, Amazon, or Microsoft have **huge data centers** that occupy **entire buildings or campuses**.

##### 💡 Example:

Google's data centers are spread across many buildings, each holding **tens of thousands of servers** running 24/7.

#### 💻 2 Server Racks and Cabinets

- The **servers** in a data center are usually kept in **rack cabinets** — tall metal frames designed to hold multiple servers stacked one above another.
- These racks are arranged in **rows**, forming **corridors or aisles** between them.
- The aisles allow technicians to access both the **front and back** of the racks easily.

#### Simple layout view:

| Rack 1 | Rack 2 | Rack 3 | → (Front Aisle)

| Rack 4 | Rack 5 | Rack 6 | → (Back Aisle)

Each rack may contain:

- Servers 
- Network switches 
- Power distribution units 
- Cooling vents 

## 5 Aisles (Corridors)

The **rows of server racks** are separated by aisles.

These are designed for:

- **Maintenance access** — technicians can walk between racks.
- **Cooling efficiency** — cold air is sent through one aisle (**cold aisle**) and hot air exits through another (**hot aisle**).

## Hot & Cold Aisle Concept:

Cold air enters from the front of the servers, and hot air comes out from the back.

Alternating the direction of racks ensures better temperature control.

## Storage Devices

Some storage devices (like large disk arrays) are **as big as entire racks**.

They are often placed **beside the server racks** for easy connection and cooling.

## Example:

A large data storage unit that holds several petabytes (PB) of information might take up one full cabinet — just like a server rack.

## 5 Large-Scale Data Centers

Modern data centers — especially those run by cloud providers — can have **thousands or even millions of servers**.

To manage so many machines efficiently, some companies use **modular or container-based data centers**.

## 6 Container-Based Data Centers

Instead of building traditional server rooms, some companies pack **1000+ servers** inside **shipping containers** (like those used on cargo ships).

Each container:

- Is **pre-built and pre-cooled**
- Can be easily **transported and installed**
- Can be **replaced as a unit** if something fails

### 💡 Example:

If a server container fails, instead of fixing each individual server, the company can **replace the entire container** — saving time and effort.

This method is called **modular data center design**.

### ⚙️ 7 Why Physical Organization is Important

Good physical organization helps in:

- **Efficient cooling** ❄️ (keeps servers from overheating)
- **Easy maintenance** 📦 (technicians can replace parts quickly)
- **Better airflow management** 🌬️
- **Space optimization** 🏙️
- **Improved reliability** ✅

## 2.4 Storage and Networking Infrastructure

### 🧠 What Does It Mean?

Every data center needs a **strong network and storage system** to connect its servers, store data safely, and communicate with users outside the data center.

In simple terms:

The **storage and networking infrastructure** of a data center is like its **nervous system** — it connects all servers together, links them to storage devices, and provides access to users on the Internet.

### 🕸️ 1 Four Main Types of Networks in a Data Center

Data centers usually use **four different kinds of physical networks**, each for a special purpose 🤝

Type of Network	Purpose	Common Technologies Used
<b>Client–Server Network</b>	Connects the data center to external users (Internet).	Ethernet or Wireless LAN
<b>Server–Server Network</b>	Enables fast communication among servers inside the data center.	Ethernet, InfiniBand
<b>Server–Storage Network</b>	Connects servers to data storage systems for reading/writing data.	Fibre Channel, Ethernet, or InfiniBand
<b>Management Network</b>	Used to monitor, control, and manage the data center equipment.	Ethernet (with separate cabling)

## ⚡ (1) Client–Server Network

This network connects the **data center** to the **outside world** (like customers, users, or web browsers).

- Example: When you open [www.amazon.com](http://www.amazon.com), your request travels from your device → Internet → Amazon's data center (through client–server network).
- Usually built using **Ethernet cables or Wi-Fi**.

### 💡 Example:

Your laptop (client) → Internet → Web server (in data center)

## ⚙️ (2) Server–Server Network

This connects **servers to each other** inside the data center.

It allows them to share workloads and data at **very high speeds**.

- Uses **Ethernet** or **InfiniBand (IBA)** for extremely fast data transfer.
- Often used in **high-performance computing (HPC)** or **cloud** environments.

### 💡 Example:

If one server stores images and another server processes them, they exchange files quickly over this network.

**Simple layout:**

Server A ↔ Server B ↔ Server C (via InfiniBand or Ethernet)

## 💻 (3) Server–Storage Network

This connects **servers** with **storage devices** (like hard drives, SSD arrays, or SANs).

- Allows servers to **store and retrieve data quickly**.
- Uses **Fibre Channel** or **Ethernet** connections for **high-speed data transfer**.
- Often managed using **Storage Area Networks (SANs)**.

### 💡 Example:

When a web app stores user data in a central database or shared storage, it uses the server–storage network.

## 🔧 (4) Management Network

- A **separate network** used by administrators to **monitor and control** the data center.
- Used for tasks like checking temperature, rebooting servers, updating configurations, etc.
- Typically uses **Ethernet**, but with **different cabling** from main networks to avoid interference.

### 💡 Example:

Admins can remotely monitor server performance or update firmware using this management network.

## 2 Load Balancer and Application Hosting

In modern data centers, **multiple applications** are hosted at the same time.

Each application has:

- Its **own servers** (physical or virtual)
- One or more **public IP addresses**

When users send requests (like visiting a website), these requests are:

1. Received at the **frontend servers**
2. Distributed evenly among multiple servers by a device called a **load balancer**

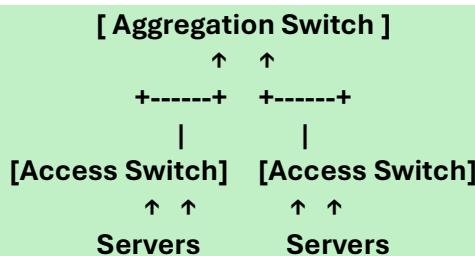
### 💡 Example:

If 100 users open a shopping site, the load balancer spreads their requests across 10 web servers — so no single server gets overloaded.

## 3 Two-Tier Network Topology (Modern Design)

Most data centers use a **two-tier network design** because it is simple and efficient.

**Structure:**



### ⚙️ How It Works

- **Access Switches:** Connect directly to servers.
- **Aggregation Switches:** Connect multiple access switches together and handle routing.

This design helps data flow efficiently between servers and users.

## 2.5 Cooling Infrastructure

### 🧠 What Is Cooling Infrastructure?

A **data center** contains **thousands of servers**, all working 24/7. These servers generate a **lot of heat** — just like how your laptop heats up after long use.

To prevent servers from overheating and failing, data centers need a **powerful and intelligent cooling system**.

So, the **cooling infrastructure** is the system that:

- Maintains the right **temperature** and **airflow** inside the data center
- Ensures **servers run efficiently** without damage
- Saves **energy** and **reduces costs**

## 1 How Cooling Works Inside a Data Center

Data centers are usually designed with **special airflow patterns** to manage heat effectively.

### The Basic Idea:

- **Cold air** is sent to the **front** of servers
- **Hot air** exits from the **back** of servers
- This air is then **cooled again and recirculated**

## 2 Hot Aisle and Cold Aisle Arrangement

To simplify cooling, **server racks** are placed in a special pattern called **Hot Aisle–Cold Aisle Design**.

**Arrangement:**

- Servers are mounted in rows (called **racks**)
- **Front sides** of two rows face each other → **Cold Aisle**
- **Back sides** of two rows face each other → **Hot Aisle**

### How It Works:

↑ Hot Air (to cooling plant)

[Rack] ← Hot Air | Cold Air → [Rack]

↑ Cold Air from raised floor

- **Cold air** is pumped up from under the floor (called a **raised plenum**) into the **cold aisles**
- **Server fans** pull cold air through the server to cool internal components
- **Hot air** comes out the back and rises into the **hot aisle**
- The **hot air** is then **sucked back to the chiller plant**, cooled again, and **recirculated**

## 3 Cooling Components Used

A typical data center cooling setup includes:

Component	Function
<b>Air Conditioning Units (CRAC/CRAH)</b>	Cools and dehumidifies air inside the data center
<b>Fans and Air Circulators</b>	Move cold and hot air through aisles
<b>Raised Floor (Plenum)</b>	Space below racks for distributing cold air
<b>Chiller Plant</b>	Cools the hot air returning from the data hall
<b>Deflectors / Containment Panels</b>	Prevent mixing of hot and cold air
<b>Temperature Sensors</b>	Monitor and maintain optimal air temperature

## Problems in Cooling Systems

Even with proper arrangements, **issues can still arise**:

- **Hot spots:** Certain servers or areas become too hot due to **uneven air distribution**
- **Air mixing:** If cold and hot air mix, cooling becomes less efficient
- **High power usage:** Cooling systems can consume **30–40% of a data center's total electricity**

### Example:

If some racks are closer to vents, they stay cool; others farther away may overheat — this creates **hot spots**.

## 5 Innovations and Green Cooling

In recent years, companies have focused on **energy-efficient (“green”) cooling systems** to reduce costs and protect the environment.

Some modern cooling methods include:

Method	Description	Example
<b>Free Cooling</b>	Uses outside air when the weather is cool	Data centers in cold regions (e.g., Finland, Canada)
<b>Liquid Cooling</b>	Uses chilled liquids or water to absorb heat	Used by Google and IBM in high-density servers
<b>Modular/Container Cooling</b>	Uses portable container units with prebuilt cooling systems	Microsoft's modular data centers
<b>Airflow Optimization</b>	Uses AI to adjust fans and airflow dynamically	Google's DeepMind AI for cooling management

## 6 Modern Data Center Designs

The concept of a **data center** has evolved based on size, purpose, and cooling methods:

Type	Description
<b>Hyperscale Facilities</b>	Huge centers (used by Google, Microsoft, Amazon) with advanced cooling and automation
<b>Consolidated Centers</b>	Merging smaller data centers into one large, efficient facility
<b>Colocation Centers</b>	Third-party facilities rented by multiple companies
<b>Hybrid Cloud Centers</b>	Combine on-premises data centers with cloud infrastructure

## 7 Site Selection and Cooling Efficiency

When choosing a location for a new data center, **cooling** and **power** play a major role.

Organizations prefer places that:

- Have **cool climates** (to save energy)
- Are **safe from natural disasters**
- Have **high-speed fiber networks**

### Examples of U.S. Data Center Hubs:

- **Silicon Valley (California)**
- **Quincy (Washington)**
- **Chicago**
- **Dallas/Fort Worth**
- **North Carolina**
- **New York/New Jersey Region**

These locations offer the **right mix of connectivity, climate, and infrastructure** for efficient cooling and operations.

## 2.6 Nature of Traffic in Data Centers

The **type of network traffic** (data flowing through the network) inside a **data center** is **very different** from traffic on the Internet.

That's because:

- Data center networks (DCNs) are **privately owned** and **optimized for high speed and low delay**.

Because of this, DCN traffic has **unique patterns and needs special handling**.

## 2 Types of Traffic in Data Centers

Traffic inside a data center is classified into **three main types** — based on size, frequency, and purpose:

### Mice Traffic, Cat Traffic, and Elephant Traffic

Let's understand each 

#### (1) Mice Traffic

Small, frequent requests sent by users — these are **short-lived** and **low-volume** data transfers.

- **Examples:**

- Google Search query
- Facebook status update
- Checking email

- **Characteristics:**

- Very short data packets
- Large in number
- Sensitive to **delay (latency)** — users expect instant responses

- **Purpose:**

Carries most **query-type traffic** in a data center.

 **Example:**

When you search on Google, your query (only a few KBs of data) travels to Google's data center — that's **mice traffic**.

 **(2) Cat Traffic**

Medium-sized messages used for **coordination** and **control** among servers.

- **Examples:**

- Downloading a small software update
- Communication between application and database servers
- Session management messages

- **Characteristics:**

- Medium data size
- **Delay-sensitive** (needs timely delivery for smooth performance)
- Occurs frequently between internal systems

- **Purpose:**

Manages **control and synchronization** within the data center.

 **Example:**

When Netflix servers coordinate to stream a video smoothly or update user data — that's **cat traffic**.

 **(3) Elephant Traffic**

Very **large data transfers** that take longer to complete.

- **Examples:**

- Uploading large backup files
- Antivirus or OS updates to all servers
- Streaming or downloading movies

- **Characteristics:**

- High volume of data
- **Throughput-sensitive** (requires high bandwidth)

- Less sensitive to delay but must maintain **steady flow**

### 💡 Example:

When YouTube updates its entire video database or when a company transfers terabytes of backup data — that's **elephant traffic**.

### 3 Table: Types of Traffic in DCNs

Type	Examples	Size / Volume	Sensitivity	Performance Requirement
<b>Mice Traffic</b>	Search queries, social media updates	Small	High delay sensitivity	Low latency
<b>Cat Traffic</b>	Control messages, small file downloads	Medium	Moderate delay sensitivity	Timely response
<b>Elephant Traffic</b>	Movie downloads, software updates, backups	Large	Low delay sensitivity, high data volume	High throughput

### 💻 ⚡ Coexistence of All Three Traffic Types

In a **data center**, all three types of traffic exist **together**:

- **Mice traffic:** Many small requests that need instant responses
- **Cat traffic:** Internal coordination among systems
- **Elephant traffic:** Bulk data transfers running in the background

This combination makes **traffic management** in DCNs very challenging — as the network must balance **speed, fairness, and efficiency** for all traffic types.

### 💡 Example Scenario (Google Data Center):

1. A user searches “weather today” → **Mice traffic**
2. Google servers communicate with databases → **Cat traffic**
3. A background process updates the language model → **Elephant traffic**

All these happen **simultaneously**, without affecting each other.

### ⚙️ 5 Why Understanding Traffic Types Matters

Knowing the nature of traffic helps engineers:

- Design **better scheduling and routing algorithms**
- Improve **load balancing**
- Reduce **network congestion**
- Enhance **overall performance and user experience**

### 3 Networking Issues in Data Centers

Networking is the **heart of cloud computing** — it connects all servers, storage systems, and users together.

Without efficient networking, cloud services like Google Drive, AWS, or Microsoft Azure wouldn't be able to function smoothly.

However, **data center networks (DCNs)** face several **issues** that affect performance, security, and reliability.

This section focuses on **three main issues**:

1.  **Availability**
2.  **Poor Network Performance**
3.  **Security**

#### **3.1 Availability**

**Availability** refers to how **accessible and operational** the cloud services are.

In simple terms, it means **how long the system stays up and running without failure**.

Cloud users expect their services to be **available 24x7** — with **zero downtime**.

#### **Problem**

Even a few seconds of service **downtime** can cause:

- Loss of **reputation**
- Violation of **Service Level Agreements (SLAs)**
- Loss of **revenue** and **customer trust**

#### **SLA (Service Level Agreement):**

A contract between a **cloud provider** and a **cloud user** that guarantees a specific **uptime** (e.g., 99.99% availability).

#### **Solution: Replication & Backup**

To ensure **high availability**, cloud providers use two main techniques:

Technique	Description	Example
<b>Replication</b>	Storing multiple copies of data in different servers or locations	Google Cloud stores your file in 3 different regions
<b>Regular Backup</b>	Taking frequent copies of data to restore in case of failure	AWS takes automatic EBS (Elastic Block Storage) snapshots

#### **Example**

Imagine you are using **Google Drive**.

If one data center fails, your files are still available from another center — this is because of **data**

**replication.**

Hence, Google maintains **high availability** even during outages.

### 3.2 Poor Network Performance

Data centers handle **massive traffic** — small, medium, and large data flows all at once (mice, cat, and elephant traffic).

Maintaining **speed and consistency** for all these types is challenging.

#### **Problem**

The **traditional TCP/IP protocol**, designed for the Internet, doesn't perform well in **data center environments** because:

- It cannot efficiently handle **bursty** and **mixed traffic**
- It introduces **latency (delay)**
- It struggles with **load balancing** among servers

#### **Key Performance Requirements of DCNs**

Requirement	Meaning	Importance
<b>High Burst Tolerance</b>	Ability to handle sudden spikes in traffic	Needed for quick queries (like search requests)
<b>Low Latency</b>	Minimal delay in data transfer	Ensures faster response to users
<b>High Throughput</b>	Ability to transfer large volumes of data efficiently	Needed for backups, video transfers, etc.

#### **Example**

In Amazon's data centers:

- Thousands of users search products (mice traffic)
- Some download product images (cat traffic)
- Others stream ads or update databases (elephant traffic)

If TCP/IP isn't optimized, small requests get delayed because large data transfers take up bandwidth — leading to **poor performance**.

#### **Solution**

To improve performance, DCNs use:

- **Custom TCP versions** (like DCTCP – Data Center TCP)
- **Load balancing** across multiple links
- **Software-defined networking (SDN)** to control and manage traffic efficiently

### 3.3 Security

Security is about **protecting data** in the cloud from loss, theft, or misuse.

Data in a data center must be protected:

- While it's **stored** (data at rest)
- While it's **being transmitted** (data in transit)

### Problems

1. **Data loss** due to power failure or backup errors
2. **Unauthorized access** or attacks by malicious insiders
3. **Physical threats** (fire, theft, hardware damage)

Even a small breach can cause **huge financial and trust damage** to the organization.

### Solutions for Security

Security Layer	Technique	Example
<b>Data Protection (in transit)</b>	Encryption during transfer	HTTPS / SSL encryption
<b>Data Protection (at rest)</b>	Encrypt stored data	AES encryption on stored files
<b>Backup &amp; Redundancy</b>	Regular automated backups	AWS S3 versioning
<b>Physical Security</b>	Restricted access, CCTV, fire safety	Biometric access control in data centers

### Example

If Microsoft Azure faces a power failure, your files aren't lost because:

- They're **encrypted**
- **Backed up** in another data center
- **Access controlled** with strict physical security

Thus, both **data safety** and **service continuity** are maintained.

### Summary Table: Networking Issues in Data Centers

Issue	Description	Example / Cause	Solution
<b>Availability</b>	Ensuring services are always online	Downtime violates SLAs	Replication, Regular Backups
<b>Poor Network Performance</b>	Slow response or congestion	TCP/IP inefficiency, bursty traffic	DCTCP, SDN, Load balancing

Issue	Description	Example / Cause	Solution
Security	Protecting data and infrastructure	Data loss, insider threats	Encryption, Backups, Physical security

## 4 Transport Layer Issues in Data Center Networks (DCNs)

The **transport layer** in networking is responsible for moving data **reliably** from one machine to another.

- **TCP (Transmission Control Protocol)** is the most widely used transport protocol on the Internet.
- TCP's **congestion control mechanisms** help avoid network overload, making it scalable and robust.

However, **data centers have unique traffic patterns** (mice, cat, and elephant traffic) and extremely **low latency requirements**, which cause some problems for TCP.

In short: **TCP works well on the Internet**, but in DCNs, its old design sometimes fails to meet all performance needs.

### 4.1 TCP Impairments in DCNs

Even after decades of evolution, TCP has **limitations** in data center networks due to:

- TCP incast
- TCP outcast
- Queue buildup
- Buffer pressure
- Pseudocongestion effect

These impairments **reduce throughput** and **increase latency**, affecting application performance.

#### 4.1.1 TCP Incast

TCP incast is a **common problem in many-to-one communication** patterns, which are frequent in DCNs.

- **Many-to-one pattern:** Multiple worker servers send data to **one aggregator server** simultaneously.
- **Problem:** All worker nodes send replies **at the same time**, causing **buffer overflow** at the **switch** connecting them.
- **Result:** Packet loss occurs, causing **delays** and **retransmissions**, especially harming **mice**

#### ⌚ Barrier-Synchronized Transmission

- In many-to-one patterns, **workers cannot send the next data block until all workers finish the current one.**

- This is called **barrier-synchronized transmission**.
- As the number of workers increases, **congestion at the switch grows**, further slowing down performance.

### 🐭 Impact on Mice Traffic

- Mice traffic (small queries) is **very delay-sensitive**.
- Frequent TCP timeouts and retransmissions **increase response time**, degrading user experience.
- **Fast retransmit** often doesn't help because **mice traffic has very few packets**, so duplicate ACKs are not enough to trigger retransmission early.

### 🛠️ Mitigating TCP Incast

Solutions exist at **different layers**:

Layer	Approach	Example
<b>Application Layer</b>	Modify how applications request or send data	Staggering worker responses
<b>Transport Layer</b>	Improve or replace TCP	DCTCP, TCP modifications
<b>Link/Network Layer</b>	Use network switches to reduce congestion	ECN marking, buffer management

### 4.1.2 TCP Outcast

TCP outcast is a **problem in data center networks** where **small flows lose out on bandwidth** when competing with **large flows** for the same output port on a switch.

- Occurs in **many-to-one patterns** (common in DCNs)
- Most **data center switches use drop-tail queues**, which worsen the problem

In simple words: **small data transfers get “pushed aside”** while large transfers dominate the network.

### ⌚ How TCP Outcast Happens

1. A **switch** has multiple input ports and one **bottleneck output port**.
2. Two sets of flows arrive:
  - **Large set of flows** (big file transfers, elephant traffic)
  - **Small set of flows** (small requests, mice traffic)
3. The **drop-tail queue** in the switch drops consecutive packets from one input port when the buffer overflows.

4. This causes **frequent timeouts** for the small flows → they get **less throughput** and higher latency.

This behavior is called a **port blackout**.

### Example Scenario (Port Blackout)

Input Port Flow Type	Result at Output Port
A      Small flows (mice)	Packets dropped consecutively → timeouts
B      Large flows (elephant)	Packets queued successfully → high throughput

- Small flows from **port A** suffer → user requests are delayed
- Large flows from **port B** dominate → take most of the bandwidth

This is counterintuitive because **TCP normally favors low RTT flows**, but TCP outcast causes **inverse RTT bias**, where **small flows with low RTT** are penalized.

### Why It Happens

1. **Drop-Tail Queues in Switches**
  - Packets are dropped **only when the buffer is full**
  - Consecutive drops → port blackout → high latency for small flows
2. **Many-to-One Communication Patterns**
  - Large and small flows arrive **at different input ports but compete for the same output port**
  - Common in DCNs for cloud applications (e.g., MapReduce, database queries)

### Mitigating TCP Outcast

To reduce TCP outcast problems, engineers use:

Method	How It Helps	Example
<b>Queue Mechanisms other than drop-tail</b>	Prevents consecutive packet drops and unfair bandwidth allocation	RED (Random Early Detection), SFQ (Stochastic Fair Queuing)
<b>End-Host Congestion Control</b>	Limits buffer occupancy at switches by controlling flow at sender	Designing efficient TCP congestion control algorithms

### 4.1.3 Queue Buildup

Queue buildup happens when **different types of traffic** (mice, cat, and elephant) share the same route in a **data center network (DCN)**.

- **Elephant traffic** (large, long-lasting transfers) can **occupy most of the network buffer**.
- **Mice traffic** (small, delay-sensitive packets) gets **delayed or dropped** as a result.

## ⚠ How Queue Buildup Affects Mice Traffic

There are **two main ways** mice traffic is affected:

### 1. Packet Loss due to Full Buffers

- Elephant traffic fills most of the switch buffers.
- Mice packets arriving at the switch may **get dropped**.
- This is similar to **TCP incast**, causing **timeouts and retransmissions**.

### 2. Increased Queuing Delay

- Even if no packets are dropped, mice packets have to **wait behind the elephant packets** in the queue.
- This increases the **response time** for delay-sensitive requests.

## 💡 Example:

- You search on Google (mice traffic) while a colleague uploads a huge backup file (elephant traffic) through the same network switch.
- Your query gets delayed because **backup packets occupy most of the buffer**, or even worse, your query packets are dropped and must be resent.

## ⚙️ Mitigating Queue Buildup

The key is to **reduce queue occupancy** in switches.

### • Problem with existing TCP:

- Most TCP variants are **reactive** → they reduce sending rate **only after packet loss**.
- By then, buffers are already congested, leading to queue buildup.

### • Desired solution:

- **Proactive congestion control** → predict congestion and reduce sending rate **before buffers overflow**.
- Helps **mice traffic get timely service** without being blocked by elephant traffic.

## 4.1.4 Buffer Pressure

Buffer pressure is a **network impairment** in data center networks (DCNs) caused when **large, long-lasting elephant traffic** monopolizes most of the **switch buffer space**.

- Mice traffic (small, delay-sensitive packets) suddenly arrives in bursts.
- Since buffers are mostly filled with elephant traffic, **mice packets get dropped**.
- This leads to **poor performance** for applications that rely on fast response times.

In simple terms: **Elephant traffic “presses down” the buffer**, leaving little room for small traffic bursts.

### ⚠ Why Buffer Pressure Happens

1. DCNs have **bursty traffic**, especially mice traffic from queries and updates.
2. Elephant traffic (large data transfers) is **greedy and long-lasting**.
3. When both coexist on the same route:
  - o Elephant traffic fills most of the buffer
  - o Mice traffic packets arriving in bursts have **nowhere to go** → packet drops
  - o This worsens latency and decreases throughput for small, important flows

### 💡 Example:

- You open a webpage (mice traffic) while a backup is being uploaded to the cloud (elephant traffic).
- The small requests for your page get **dropped or delayed** because the switch buffer is mostly full.

### ⚙️ Mitigation of Buffer Pressure

- The solution is very similar to **queue buildup mitigation**:
  - o **Minimize buffer occupancy** in switches
  - o Use **proactive congestion control** mechanisms to **prevent elephant traffic from occupying the entire buffer**
- This allows bursty mice traffic to be accommodated **without packet loss**.

### 4.1.5 Pseudocongestion Effect

Pseudocongestion occurs in **virtualized data centers** where multiple **virtual machines (VMs)** run on the same physical server.

- **TCP and UDP protocols** mistakenly perceive delays caused by **VM scheduling** as **network congestion**.
- This leads to **reduced sending rates** and **performance degradation**, even if the network is **not actually congested**.

In simple words: The network thinks it is congested, but the real problem is the **VM sharing CPU time**.

### ⚠ Why It Happens

1. **Multiple VMs on a single server**
  - o Each VM waits for the **hypervisor** to schedule CPU time.
  - o Waiting time can range from **microseconds to hundreds of milliseconds**.

## 2. TCP misinterprets delays

- High **Round-Trip Time (RTT)** due to VM scheduling → **retransmit timeout (RTO)** triggers.
- TCP reduces sending rate → throughput drops, end-to-end delay increases.

## 3. Effects observed in real-world DCNs

- Studies on **Amazon EC2** show:
  - **Throughput becomes unstable**
  - **End-to-end delay increases**, even under low network load

### Example Scenario

- Imagine a **physical server** hosting 10 VMs.
- Each VM tries to send data simultaneously.
- Hypervisor schedules each VM in turn, causing **delays**.
- TCP sender sees the delay and assumes **network congestion** → slows down transmission.

Result: The network is not truly congested, but **TCP behaves as if it is** → pseudocongestion.

### Mitigating Pseudocongestion

Two main approaches:

Approach	Description	Example
<b>Better Hypervisor Scheduling</b>	Reduce scheduling latency so VMs get CPU access faster	Efficient CPU scheduling algorithms
<b>TCP Modification</b>	Make TCP <b>intelligent</b> to detect pseudocongestion	TCP variants that differentiate <b>real congestion vs scheduling delay</b>

## 5 TCP Enhancements for DCNs

- Traditional **TCP** performs poorly in **Data Center Networks (DCNs)** due to issues such as **TCP incast, outcast, queue buildup, buffer pressure, and pseudo congestion**.
- Several **TCP variants** have been proposed to overcome these impairments.
- **three main enhancements:**
  1. **TCP with Fine-Grained RTO (FG-RTO)**
  2. **TCP with FG-RTO + Delayed ACKs Disabled**
  3. **Data Center TCP (DCTCP)**

## 5.1 TCP with Fine-Grained RTO (FG-RTO)

- **RTO (Retransmission Timeout)** defines how long TCP waits before retransmitting a lost packet.
- In standard TCP, **minimum RTO  $\approx 200$  ms**, suitable for **Internet**-scale networks where RTT  $\approx$  hundreds of ms.
- In **DCNs**, RTT  $\approx$  a few **microseconds ( $\mu$ s)**.
- Thus, a **200 ms RTO** is **too large**, causing **unnecessary waiting and throughput loss**.
- Reduce minimum RTO from **200 ms  $\rightarrow$  200  $\mu$ s**.
- This fine-grained RTO enables **faster retransmissions** and significantly improves throughput in DCNs.

### Advantages

- Requires **minimal modification** to existing TCP.
- **Easy to deploy** in existing DCN infrastructure.
- Significantly **improves throughput** by reducing delay after packet loss.

## 5.2 TCP with FG-RTO + Delayed ACKs Disabled

- **Delayed ACKs** reduce ACK overhead by sending one ACK for every **two data packets** or after a **timeout ( $\sim 40$  ms)**.
- When FG-RTO is used (RTO  $\approx 200$   $\mu$ s), this **delay (40 ms)** causes **spurious retransmissions** — sender thinks the packet is lost while receiver is still waiting to send ACK.
- Either:
  - Reduce **delayed ACK timeout** to  $\approx 200$   $\mu$ s, or
  - **Completely disable delayed ACKs** when using FG-RTO.

### Advantages

- When delayed ACK timeout is reduced or disabled, TCP achieves **higher throughput** in DCNs.
- **Avoids spurious retransmissions** that occur due to delayed ACKs.
- Works well in conjunction with **FG-RTO**.

## 5.3 Data Center TCP (DCTCP)

- Standard TCP uses **AIMD (Additive Increase / Multiplicative Decrease)** for congestion control:
  - **Additive Increase:**

$$cwnd = cwnd + 1$$

- **Multiplicative Decrease (on congestion):**

$$cwnd = cwnd/2$$

- This binary halving of cwnd is **too harsh** for DCNs, causing instability and low utilization.

### Concept of DCTCP

- **DCTCP** modifies the **multiplicative decrease** step.
- Instead of cutting cwnd by **half**, it **reduces proportionally** to the **amount of congestion** detected.
- Uses **ECN (Explicit Congestion Notification)** mechanism for fine-grained congestion feedback.

### Working Principle

1. **ECN-enabled switches** mark packets experiencing congestion.
2. **Receiver** echoes these marks back to the **sender**.
3. **Sender** estimates the **fraction of marked packets ( $\alpha$ )** — representing network congestion level.
4. The **congestion window** is reduced based on  $\alpha$ :

$$cwnd = cwnd \times \left(1 - \frac{\alpha}{2}\right)$$

- If  $\alpha = 0 \rightarrow$  no congestion  $\rightarrow$  no reduction.
- If  $\alpha = 1 \rightarrow$  full congestion  $\rightarrow$  half reduction (same as traditional TCP).

### Advantages

- Maintains **high throughput** and **low latency** simultaneously.
- **Stabilizes queue length** in switches.
- Avoids **buffer overflows** and **packet loss**, reducing incast impact.
- Proven effective in **large-scale DCNs** (e.g., **Microsoft Azure**).

#### 5.3.1 Explicit Congestion Notification (ECN)

- **ECN** is a **congestion signaling mechanism** used in IP networks to **notify** the sender about network congestion **without dropping packets**.
- It is defined in **RFC 3168** and supported in most **modern routers, switches, and operating systems**.

### How ECN Works

ECN uses **two bits** each in:

- **IP header:**
  - **ECT(0)** and **ECT(1)**  $\rightarrow$  ECN-Capable Transport
  - **CE**  $\rightarrow$  Congestion Experienced
- **TCP header:**

- **ECE** → ECN Echo (set by receiver)
- **CWR** → Congestion Window Reduced (set by sender)

### Steps in ECN Operation

1. **Negotiation Phase (3-way handshake):**
  - Sender and receiver agree to use ECN capability.
2. **Data Transmission:**
  - Sender marks packets as **ECT(0)** or **ECT(1)** to indicate ECN support.
3. **Congestion Occurs at Router:**
  - Router marks packets experiencing congestion with **CE** instead of dropping them.
4. **Receiver Action:**
  - When receiver gets a **CE-marked** packet, it sends **ACKs with ECE bit set** to inform the sender.
5. **Sender Response:**
  - After receiving ECE-marked ACKs, the sender responds with a **CWR bit** to confirm congestion handling.
  - In traditional ECN, even if the router marks **one packet**, the receiver sends **a series of ECE-marked ACKs** until the sender responds with **CWR**.
  - This ensures **reliability**, but the sender cannot determine the **amount of congestion** — only that congestion exists.

### DCTCP Modification to ECN

- **DCTCP (Data Center TCP)** modifies ECN to make it **quantitative**:
  - Receiver sends **one ECE-marked ACK only when** it receives a **CE-marked packet**.
  - This lets the sender **accurately count** the number of marked packets.
- Trade-off: reduces reliability (since if one ACK is lost, congestion info is missed).
  - But in **data centers**, ACK loss is rare → acceptable.

### DCTCP Congestion Window Update

The sender updates congestion window (**cwnd**) based on congestion fraction (**a**):

$$cwnd = cwnd \times \left(1 - \frac{\alpha}{2}\right)$$

where

$$\alpha = (1 - g) \times \alpha + g \times F$$

- **a** = fraction of packets marked by ECN ( $0 \leq a \leq 1$ )
- **F** = fraction of packets marked in previous window

- $g$  = smoothing constant ( $0 < g < 1$ )

### Switch Marking Policy

- Switch marks packets (sets **CE** bit) when its **buffer occupancy exceeds a threshold (K)**, typically **17%**.
- Sender reacts **proportionally**:
  - **Low  $\alpha$**  → **small reduction** (low congestion)
  - **High  $\alpha$**  → **large reduction** (heavy congestion)

### Advantages of DCTCP + ECN

Reduces:

- **TCP incast**,
- **queue buildup**,
- **buffer pressure**.

Provides:

- **Low latency** (for mice traffic),
- **High throughput** (for elephant traffic),
- **Burst tolerance** (for mixed workloads).

Already implemented in **Microsoft Windows Server**.

### 5.4 Incast Congestion Control for TCP (ICTCP)

- DCTCP modifies **both sender and network**, but ICTCP focuses only on **receiver-side control**.
- Goal: **Prevent packet loss due to incast by adjusting receive window (rwnd)** intelligently.
- In TCP, sender can send only up to  $\min(rwnd, cwnd)$ .
- ICTCP dynamically adjusts **rwnd** at the **receiver** to prevent congestion before it occurs.

### How ICTCP Works

#### 1. Bandwidth Quota Allocation:

- Available bandwidth is divided among all active flows as a **quota**.

#### 2. Per-Flow Control:

- Each flow's **rwnd** is increased or decreased independently.

#### 3. Adaptive rwnd Adjustment:

- Based on the **difference between expected throughput and measured throughput**:

$$\text{Adjustment Ratio} = \frac{(\text{Expected} - \text{Measured})}{\text{Expected}}$$

- If throughput is low (due to congestion), **rwnd is reduced**.

#### 4. RTT-based Estimation:

- Uses **live RTT measurements** to estimate throughput more accurately.

#### Advantages

- ✓ **No changes** needed in sender or switches — only receiver (aggregator) modified.
- ✓ Retains **backward compatibility** with standard TCP.
- ✓ Achieves **almost zero timeouts** and **high throughput**.
- ✓ Suitable for **future high-speed, low-latency networks**.

#### 5.5 – IA-TCP (Incast Avoidance TCP):

IA-TCP is designed to **overcome TCP incast problems** in data center networks (DCNs) by **controlling the rate** at which packets are sent, rather than adjusting congestion windows like traditional TCP variants (e.g., DCTCP, ICTCP).

- **Rate-based congestion control:**

IA-TCP regulates how many packets are sent into the network at once so that the **total number of outstanding packets  $\leq$  bandwidth-delay product (BDP)**.  
 → This prevents link overload and congestion collapse.

- **ACK regulation at the receiver:**

- The **receiver (aggregator)** controls the sender's rate using the **advertised window (rwnd)** field in TCP headers.
- **Minimum rwnd = 1 packet**.
- If many senders are active and the total data exceeds link capacity, IA-TCP **adds a delay ( $\Delta$ )** before sending ACKs to slow down senders and maintain network stability.

- **Desynchronization:**

The delay  $\Delta$  also helps **avoid synchronization** among senders, preventing bursty traffic.

#### Advantages:

- ✓ **Receiver-side only changes** – no need to modify senders or switches.
- ✓ **High throughput and scalability** – supports up to **96 worker nodes** in parallel without significant performance loss.
- ✓ **Improves query completion time** – critical for data-center applications.

#### 5.6 – D2TCP (Deadline-Aware Data Center TCP):

D2TCP is a **deadline-aware** TCP variant designed for **data center networks (DCNs)** that handle **high-burst and latency-sensitive** traffic. It builds upon **DCTCP** by introducing **deadline awareness** in congestion control.

- **Deadline-aware congestion control:**

D2TCP modifies the congestion window (**cwnd**) adjustment using both:

1. **ECN feedback** (to estimate congestion), and

## 2. Flow deadlines (to prioritize time-critical flows).

- **Gamma-correction function:**

This function determines how aggressively a sender reduces its cwnd based on how close the deadline is.

- **Near-deadline flows:** Back off *less* (or not at all).
- **Far-deadline flows:** Back off *more* aggressively.

- **No per-flow state:**

D2TCP maintains TCP's **distributed and reactive nature** without requiring per-flow tracking or centralized scheduling.

### Advantages:

- ✓ **Deadline awareness:** Reduces **missed deadlines by up to 75%** compared to DCTCP.
- ✓ **Built on DCTCP:** Retains DCTCP's benefits like **low queue buildup** and **incast avoidance**.
- ✓ **High burst tolerance:** Performs efficiently under heavy, bursty traffic conditions.

## 5.7 – TCP-FITDC (Fast Incast-avoiding and Throughput-enhancing Data Center TCP):

TCP-FITDC is an **adaptive delay-based TCP variant** designed for **data center networks (DCNs)** to overcome **TCP incast** and **queue buildup** problems.

It is an improvement over **DCTCP**, integrating **delay** and **ECN** information for more accurate network feedback.

- TCP-FITDC extends **DCTCP** by combining:

1. **Explicit Congestion Notification (ECN)** — to detect congestion via packet marking.
2. **Round-Trip Time (RTT) variation** — to sense queue buildup and bandwidth changes.

- It maintains **two RTT variables**:

- **rtt<sub>1</sub>:** Measured for **unmarked ACKs** (no congestion).
- **rtt<sub>2</sub>:** Measured for **marked ACKs** (congestion present).

- By comparing **rtt<sub>1</sub>** and **rtt<sub>2</sub>**, the sender estimates the **network's congestion level** more precisely.

- **Congestion window (cwnd) adjustment:**

- If no marked ACKs → Queue length below threshold → **Increase cwnd** (to boost throughput).
- If marked ACKs → Queue exceeds threshold → **Decrease cwnd** (to prevent buffer overflow).

### Advantages:

- ✓ **More accurate network estimation:** Combines ECN + RTT monitoring.
- ✓ **Better scalability** than DCTCP — up to **45 worker nodes** in parallel.
- ✓ **Prevents TCP incast** and **queue buildup**.
- ✓ **High burst tolerance** — suitable for data center traffic.

## 5.8 TDCTCP (Throughput and Delay-aware DCTCP)

TDCTCP improves upon **DCTCP** by enhancing **congestion window (cwnd) control** and **fairness**, adapting it dynamically based on congestion levels and delayed ACK behavior.

### Modifications from DCTCP

#### 1. Dynamic cwnd Increase:

Unlike DCTCP which increases cwnd steadily, TDCTCP adjusts cwnd increment based on congestion level:

$$cwnd = cwnd + \frac{1}{cwnd} + \text{function}(\alpha)$$

→ Higher increase when lightly loaded, smaller increase when heavily loaded.

#### 2. $\alpha$ Reset:

TDCTCP resets  $\alpha$  (congestion estimate) after every delayed ACK timeout to prevent stale congestion information from affecting cwnd growth.

#### 3. Dynamic Delayed ACK Timeout:

Calculates delayed ACK timeout dynamically for better fairness across flows.

### Advantages

- **26%–37% higher throughput** than DCTCP.
- **15%–20% better fairness** across flows.
- Works well in **single and multi-bottleneck topologies** and with varying buffer sizes.
- Performs better even at **low ECN marking threshold (K)**.

## 10.5.9 TCP with Guarantee Important Packets (GIP) [19]

TCP with GIP is designed to **reduce timeouts** and improve **goodput** in data center networks, especially under **many-to-one communication patterns** (common in DCNs).

## 5.10 PVTCP (Paravirtualized TCP)

PVTCP is specifically designed to address the **pseudocongestion effect** in **virtualized data centers** without requiring changes to the hypervisor.

- In virtualized environments, **hypervisor scheduling latency** can cause **variable and unpredictable RTTs**.
- This variability leads to **pseudocongestion**, where TCP mistakenly assumes network congestion and reduces sending rate unnecessarily.
- Pseudocongestion leads to **reduced throughput** and higher delays even when the network load is low.

## II. Cloud Service Providers

### 1 Introduction

Cloud computing is a technology that **provides computing resources on demand** over the Internet.

- Resources include: **software, infrastructure, platforms, and security services**.
- Unlike traditional IT, these services are **hosted by cloud providers** rather than on a company's own servers.

In simple terms: **You use computing resources as a service, without owning or managing the hardware yourself.**

#### ⚡ Key Features of Cloud Services

1. **On-demand access:** Users can access services whenever needed.
2. **Scalability:** Services can **expand or shrink** automatically to match user demand.
3. **Managed by provider:** The cloud provider handles **hardware, software, and maintenance**, so companies don't need extra IT staff.
4. **Cost efficiency:** No need to buy expensive servers or software licenses.

#### 💡 Examples of Cloud Services

- **Online storage & backup** – Google Drive, Dropbox
- **Web-based email** – Gmail, Outlook
- **Hosted office suites** – Microsoft 365, Google Workspace
- **Database processing** – Amazon RDS, Azure SQL
- **Technical support services** – Managed IT helpdesks

#### ✳️ Types of Cloud Services

Type	Description	Examples
<b>SaaS (Software as a Service)</b>	Software available over the Internet, no installation needed	Gmail, Salesforce
<b>PaaS (Platform as a Service)</b>	Provides platforms to develop and run applications	Google App Engine, Microsoft Azure App Services
<b>IaaS (Infrastructure as a Service)</b>	Provides virtualized computing resources like servers and storage	Amazon EC2, Google Compute Engine

## Popular Cloud Providers

Companies offering cloud services include:

- **Amazon, Microsoft, Google, Yahoo, EMC, Salesforce, Oracle, IBM**, and many more.

These companies provide tools and platforms to help businesses **move to the cloud efficiently** and benefit from its features.

## 2 EMC – Cloud Computing in Action

EMC is a global leader that uses cloud computing to achieve **dynamic scalability** and **infrastructure agility**, helping adapt to changing applications and business needs.

Their approach reduces complexity, optimizes infrastructure, and **cuts energy consumption** via resource sharing.

### 2.1 EMC IT – Virtualization as a Core

- **Virtualization** is the foundation of EMC IT's cloud strategy.
- It enables **resource allocation on demand**, improving **efficiency** and **utilization**.
- EMC IT delivers services in all three cloud models: **IaaS, PaaS, and SaaS**.

### Cloud Service Offerings by EMC

#### 1. IaaS (Infrastructure as a Service)

Provides infrastructure components to EMC business units:

- **Network**
- **Storage**
- **Computing**
- **Operating Systems**

These can be provisioned individually or as integrated services.

#### 2. PaaS (Platform as a Service)

Offers **secure application and information frameworks** to build and deploy solutions.

Examples of platforms provided:

- **Database Platforms**
  - Oracle Database as a Service
  - SQL Server as a Service
  - Greenplum as a Service

- **Application Platforms**

- Application Development as a Service
- Enterprise Content Management as a Service
- Information Cycle Management as a Service
- Security PaaS
- Integration as a Service

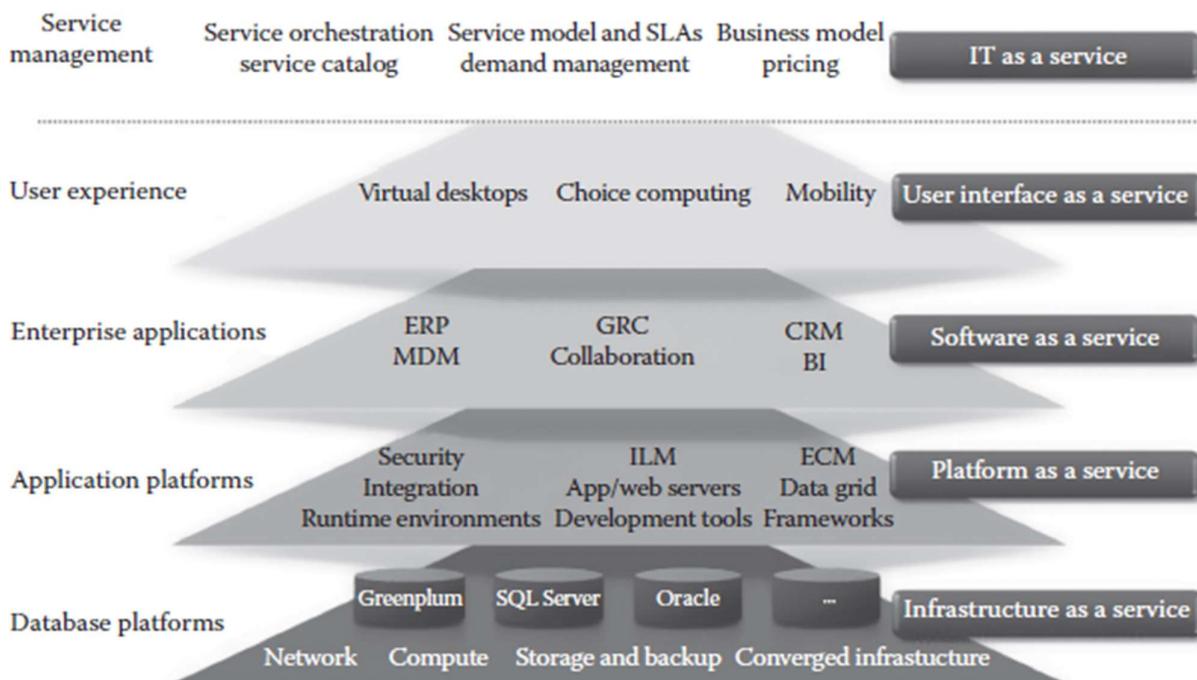
### 3. SaaS (Software as a Service)

Provides **applications and tools as services** for business enablement:

- Business Intelligence as a Service (unified architecture for multiple solutions)
- ERP (Enterprise Resource Planning) as a Service
- CRM (Customer Relationship Management) as a Service

### 4. UIaaS (User Interface as a Service)

Focuses on **provisioning user interface experiences** rather than the hardware or device itself.



## 2.2 Captiva Cloud Toolkit

EMC Captiva Cloud Toolkit is a **Software Development Kit (SDK)** designed to help developers add **scanning and imaging functionality** directly to their web-based business applications.

**Ideal for:**

- Document capture vendors
- Commercial software developers
- Enterprises creating custom, scan-enabled web applications

## Benefits:

- Developers can create a working scan-enabled application in as little as **1 week**.
- Shortens time-to-market.
- Reduces development, testing, and support costs.
- Improves return on investment.
- Accelerates competitiveness in the distributed document capture market.

## Modules of Captiva Cloud Toolkit

### 1. Scan Module

- **Function:** Imports documents from scanners.
- Works at a **page level** (imports images page by page).
- Supports multiple formats like PDF, TIFF, JPG.
- **Role:** Entry point for Captiva where all scanning starts.

### 2. MDW (Multi Directory Watch)

- **Function:** Imports documents directly from a folder/repository.
- Useful for documents arriving as soft copies (e.g., email attachments).
- Acts like a scan module but does **not** require a scanner.
- **Role:** Alternative entry point to Captiva.

### 3. IE (Image Enhancement)

- **Function:** Enhances quality of unclear images for better processing.
- Operations include:
  - Deskewing
  - Noise removal
- Configurable according to business needs.
- **Role:** Improves accuracy of data extraction.

### 4. Index Module

- **Function:** Captures key data from documents (data indexing).
- Example: For a bank form — capture account number, sort code.
- Supports validation fields to avoid incorrect data entry.
- **Role:** Central data capturing step in document processing.

### 5. Export Module

- **Function:** Sends images/data to repositories like file servers, networks, or databases.

- Used to deliver processed data to business units or departments.
- Example: Export account number and sort code to a bank's database.
- **Role:** Exit point for Captiva.

## 6. Multi Module

- **Function:** Deletes batches that have been processed and successfully exported.
- Configurable: Can be skipped if backups are needed.
- **Role:** Final cleanup step in processing.

## 3 Google

Google is one of the **leading cloud providers**, offering **secure, scalable, and reliable cloud services** to individuals and enterprises.

### Key strengths:

- Secure storage for user data
- Scalable infrastructure
- Reliable global network
- Affordable pricing — many services are free or low cost

#### 3.1 Google Cloud Platform (GCP)

Enables developers to **build, test, and deploy applications** on Google's highly advanced and scalable infrastructure.

##### 1. Highly Scalable Infrastructure

- Google has one of the largest and most advanced global networks.
- Supports rapid scaling to handle traffic spikes and demands.

##### 2. Innovative Software Infrastructure

- Includes technologies like:
  - **MapReduce** – large-scale data processing
  - **BigTable** – distributed storage system for managing structured data
  - **Dremel** – interactive analysis of large datasets

##### 3. Core Components of GCP

- **Virtual Machines (VMs)** – flexible computing power
- **Block Storage** – fast and reliable data storage
- **NoSQL Datastore** – scalable database for unstructured data

- **Big Data Analytics** – tools for large-scale data processing

#### 4. Storage Services

- Easy maintenance
- Quick access to user data

#### 5. Fully Managed & Flexible Platforms

- Managed services like **App Engine**
- Flexible VMs for custom configurations

#### 6. Automatic Scaling

- Applications can **scale up** to handle high workloads and **scale down** when traffic decreases.
- Autoscaling applies to services like App Engine and Cloud Datastore.

#### 7. Pay-as-you-go Pricing

- Users pay only for resources they consume.

### 3.2 Google Cloud Storage (GCS)

Google Cloud Storage is a **RESTful online file storage service** that allows users to store and access data over Google's infrastructure.

- **RESTful Service:** Uses REST (Representational State Transfer) architecture, which means it works via HTTP methods and is ideal for distributed systems.
- **Performance & Scalability:** Combines Google's speed and global scale.
- **Security:**
  - Data is stored redundantly across multiple physical locations for protection and reliability.
- **Sharing:** Advanced sharing capabilities make it easy to collaborate.

**Tools for Google Cloud Storage:**

#### 1. Google Developers Console

- Web-based tool for storage management tasks.

#### 2. gsutil

- Python-based command-line tool for accessing Google Cloud Storage.

### 3.3 Google Cloud Connect

Google Cloud Connect integrates **cloud storage with Microsoft Office applications** via a plug-in.

- **Cloud Master Document:**

- Saves Office documents to Google Cloud, making the cloud version the master copy.
- **Unique URL:**
  - Every document gets a unique URL for sharing.
- **Real-time Updates:**
  - Changes made by one user are visible to all who have access.
- **Conflict Resolution:**
  - If multiple users edit the same section, Cloud Connect allows selection of which changes to keep.
- **Metadata Use:**
  - Metadata identifies files so changes sync correctly.
  - Works like Google File System, leveraging Google Docs infrastructure.

#### **How It Works:**

1. Document uploaded → Metadata inserted.
2. Syncing happens → Updates sent to all copies of the document.
3. Users access the latest version with tracked changes.

#### **3.4 Google Cloud Print (GCP)**

Google Cloud Print is a **cloud-based printing service** that allows you to print documents from **any device connected to the Internet** to your printer.

#### **Requirements to Use GCP:**

1. **Google Profile** (free account).
2. **App, program, or website** that supports Google Cloud Print.
3. **Printer setup:**
  - **Cloud-ready printer** (directly connects to the Internet), or
  - Printer connected to a computer logged in to the Internet.

#### **How it Works:**

1. **Send Print Request:**
  - If using an app or website, the request is sent to Google servers.
2. **Routing:**
  - Google routes the request to the printer linked to your account.
3. **Execution:**
  - Printer executes the job if it's on, online, and has paper/ink.

#### **Printer Sharing:**

- Printers can be shared so others can send print jobs to them via Google Cloud Print.

#### Types of Printers:

##### 1. Cloud-Ready Printers:

- Directly connect to the Internet and register with Google Cloud Print.

##### 2. Non-Cloud Printers:

- Need a **connector** (via Google Chrome browser) running on a computer.
- Connector acts as a bridge between the printer and Google Cloud Print.

#### Limitations:

##### • Developer Dependency:

- Google Cloud Print needs app/website developers to integrate it into their products.
- Not every application supports it.

##### • Adoption Issues:

- Limited adoption means not every user will find the service available in all their applications.

## 3.5 Google App Engine (GAE)

Google App Engine is a **Platform as a Service (PaaS)** that lets you **build, run, and scale web applications** on Google's infrastructure — without worrying about managing servers.

#### 1. No Server Management:

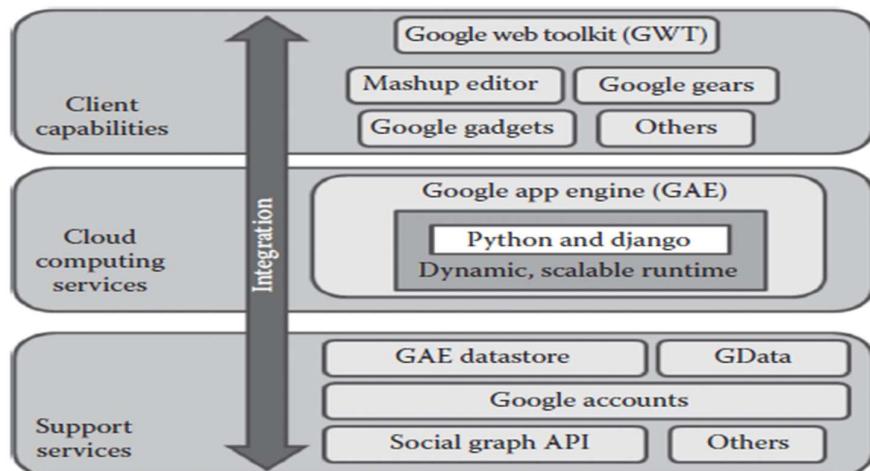
- Just upload your app — Google handles infrastructure, scaling, and maintenance.

#### 2. Domain Options:

- Serve your app on your own domain (e.g., www.example.com) or on a free subdomain (appname.appspot.com).

#### 3. Access Control:

- Apps can be **public** or restricted to organization members.



## Supported Languages & Runtime Environments:

GAE supports multiple programming environments so developers can choose freely:

Language	Environment Highlights
Java	Standard Java technologies, JVM, Java Servlets
Python	Fast interpreter, Python standard library
PHP	Native support for Google Cloud SQL & Storage
Go	Runs compiled Go code natively

These runtime environments are optimized for **speed, security, and isolation**.

## Cost Model:

- **Pay-as-you-go:** No upfront setup fees, no recurring fees.
- Pricing based on **resource usage** (storage, bandwidth, CPU) measured in GB and billed competitively.
- Control over maximum resources so your app stays within budget.

## Free Tier:

- Up to **1 GB storage** and enough CPU & bandwidth for about **5 million page views per month** — free.
- Enabling billing increases free limits, and you pay only for what you exceed.

## 4 Amazon Web Services (AWS)

AWS is a **cloud computing platform** offered by Amazon, consisting of a variety of services delivered over the Internet.

AWS provides **computing power, storage, and various cloud services** that help users avoid building their own physical infrastructure.

## Core AWS Services:

Service	Function
Amazon EC2	Virtual servers for scalable computing (IaaS)
Amazon S3	Storage service for data (object-based)
Amazon SQS	Messaging and queue service for decoupled applications

## AWS Data Centers:

Amazon has data centers across the globe, including locations in the USA (Virginia, California, Dallas, Seattle, etc.), Europe (Amsterdam, Dublin, Frankfurt, London), and Asia (Hong Kong, Singapore, Tokyo).

### 4.1 Amazon Elastic Compute Cloud (EC2)

Amazon EC2 is a **leading Infrastructure as a Service (IaaS)** that provides **scalable virtual computing capacity**.

#### 1. Virtual Machines (Instances):

- Choose from different instance types, operating systems, and software packages.
- Instances run from Amazon Machine Images (**AMI**) — pre-configured templates for EC2 instances.

#### 2. Elasticity:

- You can **increase or decrease capacity** at any time via web service interfaces.

#### 3. Interfaces:

- **AWS Management Console (GUI)** — point-and-click management.
- **Web Service API** — supports Java, PHP, Python, Ruby, Windows, .NET.

#### 4. Virtualization:

- Powered by the **Xen hypervisor**.

## EC2 Instance Types:

Type	Use Case
Standard	General purpose apps
Micro	Low throughput apps
High-memory	Memory-intensive apps
High-CPU	Compute-intensive apps
Cluster compute	High-performance computing (HPC)

## Pricing Models:

1. **On-Demand Instances:** Pay hourly without long-term commitment.
2. **Reserved Instances:** Book in advance at discounted rates.
3. **Spot Instances:** Bid on unused capacity to get cheaper rates (price fluctuates).

## Deployment Locations:

- **Regions:** Geographical areas with multiple data centers.
- **Availability Zones (AZs):** Isolated locations within regions to improve fault tolerance.

#### **Management & Monitoring:**

- **Amazon CloudWatch:** Monitors CPU, disk I/O, and network usage.
- **AWS Management Console** and APIs allow full control over EC2 instances.

#### **Security:**

- Instances use a **signature-based authentication** via key pairs.
- **Amazon VPC** allows secure connection of existing infrastructure via VPN.

#### **Storage Options:**

##### **1. Elastic Block Storage (EBS):**

- Persistent block storage volumes that can be attached to EC2 instances.
- Can boot an EC2 instance or be used as additional storage.

##### **2. Amazon Simple Storage Service (S3):**

- Highly durable object storage for mission-critical data.
- Stores **objects** in **buckets**, accessed via unique developer-assigned keys.

#### **Elastic Load Balancing (ELB):**

- Automatically distributes incoming traffic among instances.
- Provides fault tolerance by routing away from unhealthy instances.

### **4.2 Amazon Simple Storage Service (Amazon S3)**

Amazon S3 is a **highly scalable, secure, and durable cloud storage service** designed to make storage easy for developers and organizations.

Amazon S3 is essentially "**storage for the Internet**", enabling users to store and retrieve any amount of data at any time, from anywhere in the world.

#### **Features of Amazon S3:**

<b>Feature</b>	<b>Description</b>
<b>Simple Web Interface</b>	Allows storing and retrieving data via a web service interface.
<b>Scalability</b>	Automatically scales to handle any amount of data without manual intervention.
<b>Durability &amp; Availability</b>	Stores data redundantly across multiple devices and locations for high durability.
<b>Security</b>	Offers encryption, access control, and secure data transfer.
<b>Performance</b>	Fast access to stored data anywhere in the world.
<b>Cost Efficiency</b>	Pay-as-you-go pricing model.

## How Amazon S3 Works:

- **Objects:** Data in S3 is stored as **objects** (files).
- **Buckets:** Objects are stored in **buckets** (containers).
- **Keys:** Each object has a unique key (name) used to retrieve it.

## Special Features of Amazon S3:

1. **Reduced Redundancy Storage (RRS):**
  - Stores data at lower redundancy for lower cost.
  - Useful for data that can be reproduced or stored elsewhere.
2. **Integration with EC2:**
  - Data stored in S3 can be used by Amazon EC2 without data transfer charges between services.
3. **Data Analytics:**
  - S3 is ideal for storing large datasets for analytics, such as pharmaceutical or financial data.
4. **Backup & Archiving:**
  - Highly durable storage for critical backups and archives.
5. **AWS Import/Export:**
  - Transfer large amounts of data into or out of AWS using physical storage devices.
  - Useful for periodic backups or disaster recovery.
6. **Static Website Hosting:**
  - Host static content like HTML, images, videos, and JavaScript directly from S3.

## Use Cases:

- **Web application storage** — store assets like images, videos, and HTML files.
- **Data analytics storage** — store large datasets for computation.
- **Backup & disaster recovery** — highly durable storage for critical data.
- **Website hosting** — hosting static websites.

## 4.3 Amazon Simple Queue Service (Amazon SQS)

Amazon SQS is a **fully managed, reliable, scalable message queuing service** provided by AWS.

It allows **decoupling of application components**, meaning different parts of an application can communicate without being directly connected, improving reliability and scalability.

Amazon SQS acts as a **temporary message repository (queue)** where messages are stored until they are processed by other components of the application.

Example:

- Application A produces messages → Sends to SQS queue
- Application B consumes messages → Retrieves from SQS queue

This ensures **smooth communication without losing data** even if one component is busy or temporarily unavailable.

### **Features of Amazon SQS:**

<b>Feature</b>	<b>Description</b>
<b>Fully Managed</b>	No need to manage servers or infrastructure.
<b>Reliable</b>	Guarantees message delivery through redundant storage across multiple servers and data centers.
<b>Scalable</b>	Can handle any volume of messages and throughput.
<b>Decoupling Components</b>	Components work independently without direct dependencies.
<b>Multiple Readers &amp; Writers</b>	Supports multiple producers and consumers simultaneously.
<b>Variable Length Messages</b>	Supports different sizes of messages.
<b>Access Control</b>	Fine-grained permissions for each queue.
<b>Configurable Queues</b>	Custom settings such as message retention time and visibility timeout.

### **How Amazon SQS Works:**

- 1. Message Creation:**
  - A producer sends a message to the queue.
- 2. Message Storage:**
  - The message is stored in the queue until processed.
- 3. Message Processing:**
  - A consumer retrieves the message and processes it.
- 4. Message Deletion:**
  - After successful processing, the consumer deletes the message from the queue.

## Types of Queues in Amazon SQS:

### 1. Standard Queues

- High throughput.
- Messages might be delivered more than once.
- Order of messages not guaranteed.

### 2. FIFO (First-In-First-Out) Queues

- Exactly-once message processing.
- Strict order of message delivery.

## Use Cases:

- **Order processing systems** — decoupling order intake and order processing.
- **Microservices communication** — connecting services without tight integration.
- **Batch processing** — storing jobs to be processed later.
- **Scaling workloads** — smoothing traffic spikes by queuing messages.

## 5 Microsoft

### Microsoft and Cloud Computing (MSIT)

Microsoft views cloud computing as the **preferred default environment** for new and migrated applications.

Microsoft IT (MSIT) uses a **methodology and best practices** to analyze which applications are suitable for migration to the cloud.

This enables MSIT to **choose the ideal cloud computing environment** for each application and helps customers migrate effectively.

### 5.1 Windows Azure (Microsoft Azure)

**Windows Azure** is Microsoft's cloud computing platform providing **Platform as a Service (PaaS)** and **Infrastructure as a Service (IaaS)** solutions.

Key offerings:

- **Web Roles** — for hosting front-end web applications.
- **Worker Roles** — for background processing tasks.

**Management:** Azure handles operating systems and server management so developers can focus on application development.

### Tools related to Windows Azure:

#### 1. Migration Assessment Tool (MAT)

- Helps assess readiness for migration to Azure.
- Produces reports on development effort and architectural needs.

## 2. Windows Azure Pricing Calculator

- Compares current infrastructure costs with Azure cloud costs.

## 3. Windows Azure Pack for Windows Server

- Lets customers install Azure technologies in their own data centers.
- Runs on Windows Server 2012 R2 and System Center 2012 R2.
- Offers a **self-service, multi-tenant cloud** consistent with Azure public cloud experience.

## 5.2 Microsoft Assessment and Planning Toolkit (MAP)

**MAP Toolkit** is a **free, agentless planning and assessment tool** designed for cloud migration readiness.

Key features:

- Provides detailed **readiness assessment reports**.
- Generates **executive proposals** for migration.
- Collects hardware and software inventory.
- Gives **recommendations** to speed up migration to public or private clouds.
- Analyzes **server utilization** for virtualization and consolidation with **Hyper-V**.

## 5.3 SharePoint

**Microsoft SharePoint** is a **web-based collaboration and document management platform**.

Key features:

- Intranet portals, document/file management, and collaboration tools.
- Enterprise social networking and extranets.
- Websites and enterprise search functionality.
- Business intelligence, system integration, and workflow automation.
- Microsoft Office-like interface for ease of use.

**Difference from Google Cloud Connect:**

- SharePoint is **not free**.
- Offers more advanced collaboration features that go beyond Google's offerings.

## 6 IBM and Cloud Computing

IBM is a **major cloud computing player** offering a variety of cloud services under the brand name **IBM SmartCloud**.

IBM's cloud services include **Infrastructure as a Service (IaaS)**, **Platform as a Service (PaaS)**, and **Software as a Service (SaaS)** through **public, private, and hybrid cloud models**.

IBM also provides hardware, middleware, and management tools to build, manage, and scale cloud solutions for enterprises.

## 6.1 IBM Cloud Models

IBM offers flexible cloud deployment models ranging from fully **private cloud** to fully **public cloud**, and combinations in between (hybrid cloud).

### Five Cloud Provision Models from IBM:

1. **Private cloud** — Owned and operated by the customer.
2. **Private cloud** — Owned by the customer but operated by IBM or another provider.
3. **Private cloud** — Owned and operated by IBM or another provider.
4. **Virtual private cloud services** — Multi-tenant support for individual enterprises.
5. **Public cloud services** — Functions provided to individuals over the Internet.

### Specialized Private Cloud Solutions from IBM:

- **IBM Workload Deployer:** Connects existing servers to virtualization components and middleware for easy cloud deployment.
- **IBM Cloudburst:** A "cloud-in-a-box" solution containing blade servers, middleware, and virtualization tools to build private clouds.

IBM also offers **hardware and software building blocks** and **reference architectures** for customers who want to integrate and deploy their own private clouds.

## 6.2 IBM SmartCloud

**IBM SmartCloud** is IBM's ecosystem for cloud computing, combining hardware, software, and services under a single framework.

IBM organizes SmartCloud into three main parts:

### 1. SmartCloud Foundation

- Core infrastructure for cloud deployment.
- Includes hardware, provisioning, management, integration, and security tools.
- Serves as the foundation for both **private and hybrid clouds**.

### 2. SmartCloud Services

- Includes **IaaS**, **PaaS**, and backup services.
- Examples:
  - Infrastructure services (storage, computing, networking resources).
  - Platform services (development frameworks, runtime environments).
  - Application services (SaaS).

### 3. SmartCloud Solutions

- SaaS applications designed for collaboration, analytics, and business marketing.
- Includes **Business Process as a Service (BPaaS)**: Business processes delivered through the cloud with **self-service provisioning, elastic scaling, and usage metering**.

## 7 SAP Labs

**SAP Labs** develops enterprise software for managing business operations and customer relations.

SAP is a **global leader in enterprise applications** with products in ERP (Enterprise Resource Planning), business analytics, mobile computing, and in-memory computing.

**Key products include:**

- **SAP ERP** – Enterprise Resource Planning systems
- **SAP BW (Business Warehouse)** – Data warehousing solution
- **SAP Business Objects** – Business intelligence software
- **Sybase mobile products**
- **SAP HANA** – In-memory computing platform

SAP is one of the largest software companies globally, offering powerful enterprise solutions.

### 7.1 SAP HANA Cloud Platform (HCP)

**SAP HANA Cloud Platform** is an **open-standard, modular Platform as a Service (PaaS)**.

It is **Eclipse-based**, meaning developers can build and deploy applications using familiar tools.

**How it works:**

- Applications are deployed to the cloud as:
  - **WAR files** (Web Application Archive)
  - **OSGi bundles** (Java jar components with extra manifest headers)
- Runs inside a **Java-based SAP HANA Cloud Platform runtime environment**
- Managed via **web-based tools**

**Main Features of SAP HANA Cloud Platform:**

- **Enterprise platform for developers**: Quick application building and deployment
- **Native integration**: Works seamlessly with SAP and non-SAP systems
- **In-memory persistence**: Real-time processing and data analysis
- **Secure data platform**: Built with strong security features
- **Lightweight modular runtime**: Efficient and scalable application environment

## 7.2 Virtualization Services Provided by SAP

**ERP Virtualization** offers benefits such as:

- **Shorter development cycles**
- **Reduced IT costs**
- **Improved system availability**
- **Energy savings**

SAP partners with **VMware** to offer services that transition enterprise applications to a **private cloud platform** based on proven virtualization technology.

**Business benefits of ERP virtualization:**

- Maximizes hardware utilization
- Improves efficiency and resource management
- Enables flexible cloud-based ERP deployment

## 8 Salesforce

**Salesforce.com** is a **cloud computing and social enterprise SaaS provider** headquartered in San Francisco.

It is best known for **Salesforce CRM** (Customer Relationship Management) and a suite of related services and platforms.

**Offerings:**

- **Sales Cloud** – CRM for sales processes
- **Service Cloud** – CRM for customer service
- **Marketing Cloud** – Digital marketing tools
- **Force.com** – Platform for building custom applications
- **Chatter** – Social collaboration tool
- **Work.com** – Performance management platform

**Other key offering:**

- **AppExchange** – Marketplace for custom applications, integrations, and tools
- Consulting, deployment, and training services

### 8.1 Sales Cloud

**Purpose:** A complete sales module designed to manage the entire sales process from lead generation to closing deals.

**Core Components:**

- Leads
- Accounts

- Contacts
- Contracts
- Opportunities
- Products
- Pricebooks
- Quotes
- Campaigns

**Features:**

- **Web-to-lead** capture and autoresponse
- Centralized contact and deal management
- Integration with social media
- Real-time updates pushed automatically (contact info, deal updates, approvals)
- Easy-to-use, cloud-based platform
- Open architecture with automatic updates, eliminating traditional CRM complexities

**Benefit:**

Enables a smooth, end-to-end sales process that is accessible anytime, anywhere, without expensive installations or upgrades.

## 8.2 Service Cloud: Knowledge as a Service

**Purpose:** A customer service module that focuses on delivering efficient and intelligent customer support.

**Core Components:**

- Accounts
- Contacts
- Cases
- Solutions

**Features:**

- **Public Knowledge Base** – Centralized repository of information
- **Web-to-case** – Create cases directly from the web
- Call center support
- Self-service portals
- Customer service automation
- Social networking integration (e.g., Twitter, Facebook)

### How It Works:

- Service Cloud hosts the first **enterprise-grade knowledge base** on a **multitenant cloud platform**.
- Knowledge becomes a **process** that can be continually created, reviewed, delivered, analyzed, and improved.
- Agents can instantly access information via phone, email, or chat.
- Knowledge can be shared publicly or across social platforms for better customer engagement.
- Fully integrated with CRM for seamless service operations.

### Benefit:

Provides powerful customer service tools without the need for costly, on-premises infrastructure, while keeping knowledge accessible and secure.

## 9 Rackspace

**Rackspace Cloud** is a prominent **Infrastructure-as-a-Service (IaaS)** provider offering reliable, scalable cloud computing solutions.

It provides three core solutions:

1. **Cloud Servers** – On-demand computational power
2. **Cloud Files** – Elastic online file storage and content delivery
3. **Cloud Sites** – Robust, scalable web hosting

### 9.1 Cloud Servers

Cloud Servers are **virtual machines (VMs)** that run in Rackspace's cloud, providing **computing capacity on demand**.

- **Flavors** – Hardware configurations (disk space, memory, CPU priority).
- **Images** – Pre-built or custom system templates for creating servers.

Examples:

- Linux: Ubuntu, Debian, Gentoo, CentOS, Fedora, Arch, Red Hat Enterprise Linux
- Windows: Windows Server 2008, Windows Server 2003

### Virtualization

- Linux: Xen Hypervisor
- Windows: Xen Server

### Server Sizes

- Memory: 256 MB to 15.5 GB
- CPU power: Extra CPU provided automatically when available, free of cost.

## 9.2 Features & Operations

### Backup & Customization

- Backup schedules to create server images for recovery
- **Gold server images** for frequently used configurations

### Management

- **Rackspace Cloud Control Panel (GUI)** – Billing, reporting, support resources
- **Cloud Server API (RESTful)** – Programmatic control, open-sourced under Creative Commons Attribution 3.0
- Supported bindings: C++, Java, Python, Ruby

### Security

- Token-based authentication using HTTP x-Header
- Public/private keys for secured shell (SSH) access

## 9.3 Scalability & Load Balancing

- **Auto-scaling** – Load balancing initiated via Control Panel or API
- Adjusts RAM, disk space, CPU allocation
- Server temporarily taken offline during scaling
- Rackspace offers **Cloud Load Balancing** (beta) as a dedicated service
- Servers can be configured as load balancers using open-source tools

## 9.4 Storage

- **Persistent storage** through RAID10 disk storage
- Guarantees data persistence and reliability

## 10 VMware

VMware is a **leader in virtualization technology** and offers **enterprise cloud computing solutions for private, public, and hybrid clouds**.

VMware's core offerings include:

- **VMware vSphere**
- **VMware vCloud Director**
- **VMware vShield**
- **VMware vCloud Datacenter Services**
- **VMware vCloud Express**

## 10.1 Private Cloud with VMware

**Private clouds** allow better utilization and management of internal IT infrastructure compared to traditional systems.

Benefits include:

- Greater operational efficiency
- Enhanced security
- Fault tolerance
- Standardization
- Rapid provisioning
- Self-service for applications
- Cost savings through infrastructure consolidation

### VMware Products for Private Cloud

- **VMware vSphere**
  - Virtualization platform to transform IT infrastructure into virtual compute, storage, and network resources
  - Provides management at both infrastructure and application levels
  - Enables efficient operations and scalable, secure applications
- **VMware vCloud Director**
  - Works with vSphere to build **secure, multitenant private clouds**
  - Pools infrastructure into **virtual datacenters**
  - Provides **web-based portals** and **programmatic interfaces** for automated, catalog-based service delivery
  - Features:
    - Isolated virtual resources
    - LDAP authentication
    - Policy controls
    - Unique catalogs
- **VMware vShield**
  - Security for cloud environments
  - Services include:
    - Perimeter protection
    - Port-level firewall

- NAT and DHCP
- Site-to-site VPN
- Network isolation
- Web load balancing

## 10.2 Public & Hybrid Clouds

VMware provides **public and hybrid cloud solutions** via partnerships with certified service providers.

### Key Offerings

- **VMware vCloud Datacenter Services**
  - Extends private cloud with external resources
  - Scalable environment built on vCloud Director and vSphere technology
  - Enables interoperability between clouds
  - Allows bursting of private clouds into public clouds
- **VMware vCloud Express**
  - IaaS offering from VMware service provider partners
  - On-demand, pay-as-you-go infrastructure
  - Providers: Virtacore, Hosting.com, Melbourne IT, Terremark's vCloud Express
  - Flexible instance types, load balancing, storage, and pricing

## 10.3 Features

Feature	Details
<b>Private Cloud</b>	Standardized, secure, scalable, cost-efficient
<b>Virtualization</b>	Managed by VMware vSphere
<b>Multitenancy</b>	Enabled via vCloud Director
<b>Security</b>	vShield services
<b>Hybrid Cloud</b>	vCloud Datacenter Services
<b>On-demand Public Cloud</b>	vCloud Express