# CLOUD COMPUTING

## UNIT I
## FUNDAMENTAL CLOUD COMPUTING AND VIRTUALIZATION

### 1.0 Origin and influences

The idea of computing in a "cloud" traces back to the origins of utility computing, a concept that computer scientist John McCarthy publicly proposed in 1961:

*"If computers of the kind I have advocated become the computers of the future, then computing may someday be organized as a public utility just as the telephone system is a public utility. … The computer utility could become the basis of a new and important industry."*

In 1969, Leonard Kleinrock, a chief scientist of the Advanced Research Projects Agency Network or ARPANET project that seeded the Internet, stated:

*"As of now, computer networks are still in their infancy, but as they grow up and become sophisticated, we will probably see the spread of 'computer utilities' …".*

The general public has been leveraging forms of Internet-based computer utilities since the mid-1990s through various incarnations of search engines (Yahoo!, Google), e-mail services (Hotmail, Gmail), open publishing platforms (MySpace, Facebook, YouTube), and other types of social media (Twitter, LinkedIn).

In the early 1990s, Salesforce.com pioneered the idea ofbringing remote configuration services to the enterprise. In 2002, Amazon.com launched the Amazon Web Services (AWS) platform, a set of businessfocused services that provide storage, computing resources, and business operations.

In the early 1990s, slightly different versions of the term "network cloud" or "cloud" were introduced in the Internet industry. Although the word "cloud" is also used for mobile phones, it refers to the abstraction process of data transfer methods between public and semi-public networks, mainly packet switching.

It wasn't until 2006 that the term "cloud computing" emerged in the business world. It was around this time that Amazon launched its Elastic Compute Cloud (EC2) service, which allows organizations to "rent" computing power and processing power to run business applications.

### 1.0.1 Definitions

**By Gartner (2008):**

*"…a style of computing in which scalable and elastic IT-enabled capabilities is delivered as a service to external customers using Internet technologies."*

**By Forrester Research**

*"…a standardized IT capability (services, software, or infrastructure) delivered via Internet technologies in a pay-per-use, self-service way."*

**By NIST (2011):**

*"Cloud computing is a model for enabling ubiquitous, convenient, on-demand network access to a shared pool of confi gurable computing resources (e.g., networks, servers, storage, applications, and services) that can be rapidly provisioned and released with minimal management effort or service provider interaction. This cloud model is composed of fi ve essential characteristics, three service models, and four deployment models."*

### 1.0.2 Business Drivers

- Capacity Planning

- Cost reduction

- Organizational Agility

### 1.0.3 Technology Innovations

- Clusteriing

- Grid Computing

- Virtualization

### 1.0.4 Cloud Enabling Technologies

- Broadband Networks and Internet Architecture

- Data Center Technology

- (Modern) Virtualization Technology

- Web Technology

- Multitenant Technology

- Service Technology

## 1.1 Data access and integration

Data access and integration are critical aspects of leveraging cloud computing effectively. They involve ensuring seamless access to data stored in the cloud and integrating this data with existing applications and workflows. Here's a detailed exploration of data access and integration in the context of cloud computing:

### 1.1.1 Data Access in the Cloud

1. **Access Methods**
   - **APIs**: Cloud providers offer APIs (Application Programming Interfaces) that allow applications to interact programmatically with cloud services. APIs facilitate actions such as reading, writing, and managing data stored in the cloud.
   - **Web Interfaces**: Many cloud services provide web-based portals or consoles where users can log in and manage their data directly through a user-friendly interface.

- o **Command-Line Interfaces (CLI)**: Some cloud platforms offer CLI tools that allow developers and administrators to manage cloud resources via command-line commands.
2. **Security and Authentication**
    - o **Identity and Access Management (IAM)**: Implementing IAM policies to control who can access data and resources within the cloud environment.
    - o **Encryption**: Encrypting data both in transit and at rest to protect it from unauthorized access.
    - o **Access Controls**: Applying granular access controls to ensure that only authorized users and applications have access to specific data.
3. **Data Transfer**
    - o **Data Transfer Methods**: Using secure protocols such as HTTPS, FTPS, or SCP to transfer data between on-premises systems and the cloud.
    - o **Data Migration Services**: Cloud providers often offer data migration services that facilitate the transfer of large volumes of data into and out of the cloud.

## 1.1.2 Data Integration in the Cloud

1. **Challenges**
    - o **Data Formats**: Ensuring compatibility between different data formats used by cloud services and existing on-premises systems.
    - o **API Compatibility**: Ensuring that APIs used by cloud services are compatible with the APIs and protocols used by existing applications.
    - o **Latency**: Minimizing latency when accessing and transferring data between cloud services and on-premises systems.
    - o **Data Consistency**: Maintaining data consistency and synchronization across distributed systems.
2. **Strategies**
    - o **APIs and Middleware**: Using APIs and middleware to facilitate communication and data exchange between cloud services and existing applications.
    - o **Data Replication and Synchronization**: Implementing mechanisms to replicate and synchronize data between cloud and on-premises systems to ensure consistency.
    - o **Event-Driven Architecture**: Adopting event-driven architectures where changes in data trigger events that propagate across integrated systems.
    - o **Data Virtualization**: Using data virtualization techniques to create a unified view of data across disparate sources, including cloud and on-premises environments.
3. **Integration Patterns**
    - o **Batch Processing**: Moving large volumes of data in scheduled batches between systems.
    - o **Real-Time Integration**: Enabling real-time data exchange and processing between cloud services and existing applications.
    - o **Streaming Data**: Handling continuous streams of data generated by IoT devices, sensors, or real-time analytics applications.

## 1.1.3 Best Practices

1. **Plan and Architect**: Design a comprehensive data access and integration strategy that aligns with business objectives and IT infrastructure.
2. **Security**: Prioritize data security by implementing encryption, access controls, and monitoring to protect sensitive data.
3. **Scalability**: Ensure that data access and integration solutions can scale to accommodate growing volumes of data and increasing user demand.
4. **Monitoring and Optimization**: Continuously monitor data access patterns and integration processes to identify bottlenecks and optimize performance.

## 1.2 Basic concepts and terminology

- **Cloud**

  Cloud refers to a unique IT environment designed to provide scalable and s calable IT resources. The term is actually a metaphor for the Internet, a ne twork that provides remote access to IT systems. Before cloud computing b ecame the official IT business area, cloud symbols were often used to repre sent the Internet as a web architecture in many private and important docu ments.
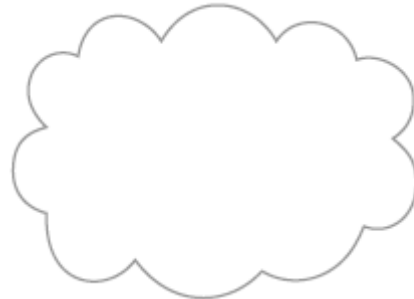
**Fig: Symbol for cloud**

- **IT resources**

  IT **resources are** physical or virtual IT-related **assets** that can be **software-based (such** as virtual **servers** or **software)** or **hardware-based (such** as physical **servers** or network **equipment).**

physical server    virtual server    software program    service    storage device    network device
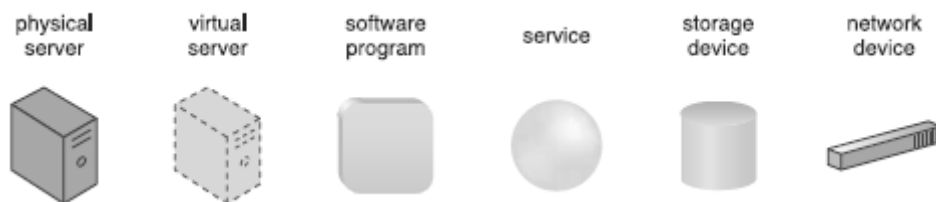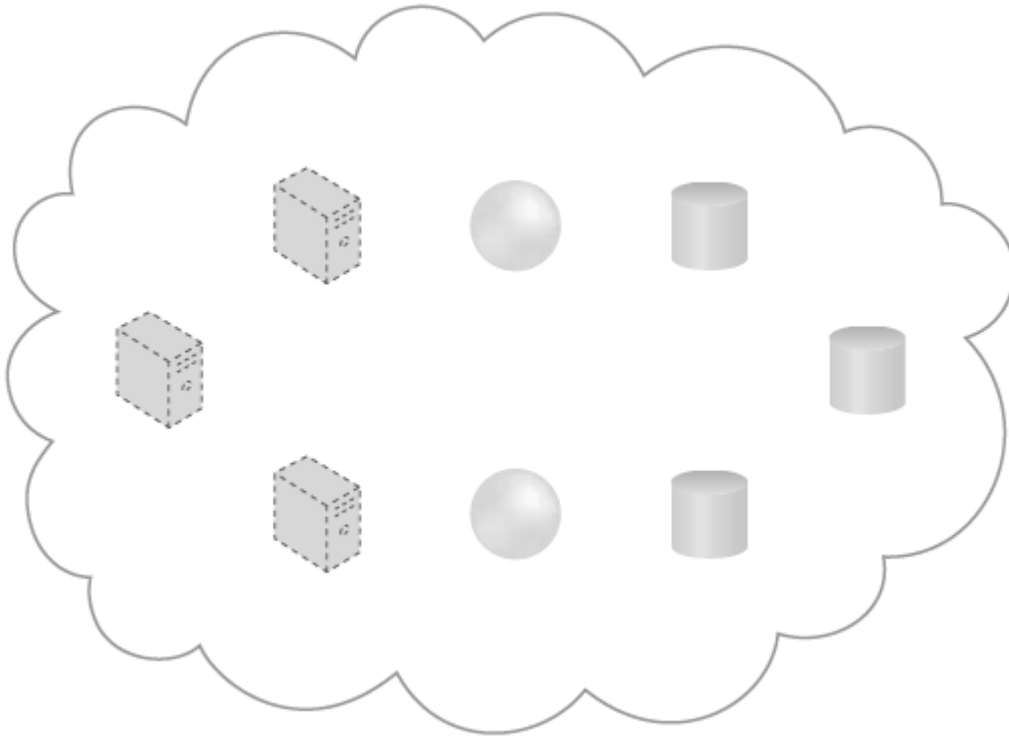
**Fig: IT Resources**

**Fig: Cloud hosting IT resources**

- **On premise**
  As a unique and accessible remote location, the cloud represents an option for outsourcing IT resources. IT resources in the IT business are always at the edge (organizations not specifically represented in the cloud), considered to be in the IT business domain or solely in the domain domain. In other words, the term "onsite" is just another way of saying "in the context of managing a non-cloud-based IT environment."

- **Cloud Consumers and Cloud Providers**
  The cloud provider is the party that provides cloud-based IT services. Parties using cloud IT resources are cloud users. These terms represent the responsibilities that organizations typically assume regarding the cloud and its contractual relationship.
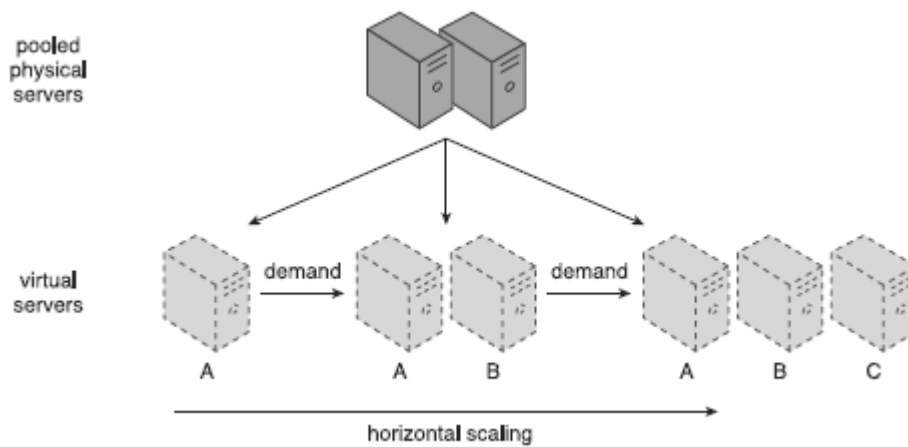- **Scaling**

**Horizontal & vertical Scaling**
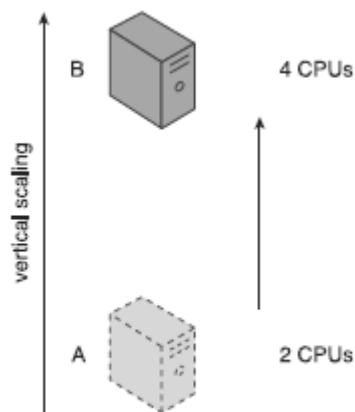


Fig: Horizontal scaling



Fig: Vertical scaling

| Horizontal Scaling | Vertical Scaling |
|---|---|
| less expensive (through commodity hardware components) | more expensive (specialized servers) |
| IT resources instantly available | IT resources normally instantly available |
| resource replication and automated scaling | additional setup is normally needed |
| additional IT resources needed | no additional IT resources needed |
| not limited by hardware capacity | limited by maximum hardware capacity |

- **Cloud Service**
  Cloud services are all IT services that can be accessed remotely from the cloud. Unlike other areas of IT that fall under the service technology umbrella, such as service-
  oriented architecture, the term "service" is particularly broad in the context of cloud counting. Cloud services can be simple web-
  based software programs with communication tools called messaging, as well as managed devices or remote access to larger areas and other IT services.
- **Cloud Service Consumer**
  Cloud service users are responsible for the temporary uptime assumed by the software program when accessing the cloud service. Types of cloud services may include offices that can access cloud services and other IT areas such as mobile devices, manuals work, as well as software services and services that can be accessed through the cloud service by arranging contracts.



**Fig: Cloud Service consumers**

### 1.2.1 Goals & benefit
- Reduced Investments and Proportional Costs
- Increased Scalability
- Increased Availability and Reliability

### 1.2.2 Risk & Challenges
- Increased Security Vulnerabilities
- Reduced Operational Governance Control
- Limited Portability Between Cloud Providers
- Multi-Regional Compliance and Legal Issues

## 1.3 Cloud Deployment Models

The National Institute of Standards and Technology (NIST) identify four primary types of cloud deployment models, each catering to different organizational needs and scenarios. Here's an overview of each:
- ➢ Private cloud
- ➢ Community cloud
- ➢ Public cloud
- ➢ Hybrid cloud

### 1.3.1 Private Cloud

Cloud environments can be categorized based on hardware location and ownership. Here's a summary of what you described:

- **Accessibility**: Only accessible to a specific organization.
- **Location**: Can be on-site (on the organization's premises) or off-site (hosted by a third party).
- **Ownership and Management**: Operated solely for one organization. Management can be handled internally by the organization's IT department or outsourced to a third-party provider.
- **Security and Control**: Typically offers greater security and control compared to public clouds because it is dedicated to a single organization.

**Key Points**

- A private cloud provides the benefits of cloud computing (such as scalability and resource efficiency) while maintaining a higher level of security and control.
- The location of the private cloud can vary; it does not have to be on the organization's premises to qualify as a private cloud.
- Management of the private cloud can be done by the organization itself or by a third-party provider.

This makes private clouds a flexible and secure option for organizations that need to comply with strict security and regulatory requirements while still benefiting from cloud technology.
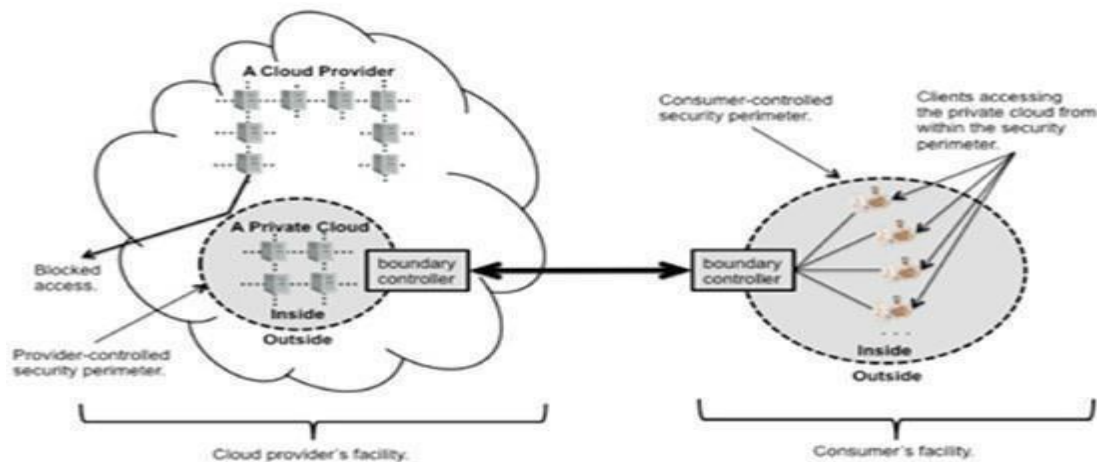


**Figure 1.1 Schematic Sketch of Private Cloud**

### 1.3.2 Community Cloud
Community cloud, which is another type of cloud deployment model. Here's a detailed overview:
- **Community Cloud**
- **Accessibility**: Shared by several organizations that have common concerns or objectives (e.g., mission, security requirements, policy, and compliance considerations).
- **Location**: Can be on-site, off-site, or a combination of both, depending on the agreements among the participating organizations.
- **Ownership and Management**: The cloud infrastructure is collaboratively managed and used by the participating organizations. It can be managed by the organizations themselves or by a third-party service provider.
- **Security and Control**: Offers a middle ground between private and public clouds in terms of security and control. While it is more secure than a public cloud because it is not open to the general public, it may not provide the same level of control as a private cloud since it is shared among multiple organizations.
- **Key Points**
- **Shared Resources**: Resources and infrastructure are shared among the participating organizations, leading to cost savings and efficient resource utilization.
- **Common Interests**: Typically used by organizations with shared concerns or interests, such as government departments, universities, central banks, and other institutions with similar requirements.
- **Compliance and Security**: Ensures compliance with specific

regulations and security policies that are common to the participating organizations.

- **Use Cases**
- **Government Departments**: Different departments within a government can share a community cloud to streamline operations and ensure compliance with governmental regulations.
- **Universities**: Academic institutions can collaborate and share resources for research and educational purposes.
- **Central Banks**: Financial institutions with strict security and compliance requirements can benefit from a community cloud tailored to their needs.
- Community clouds provide a balance of cost-efficiency, security, and control, making them suitable for organizations with shared goals and requirements.
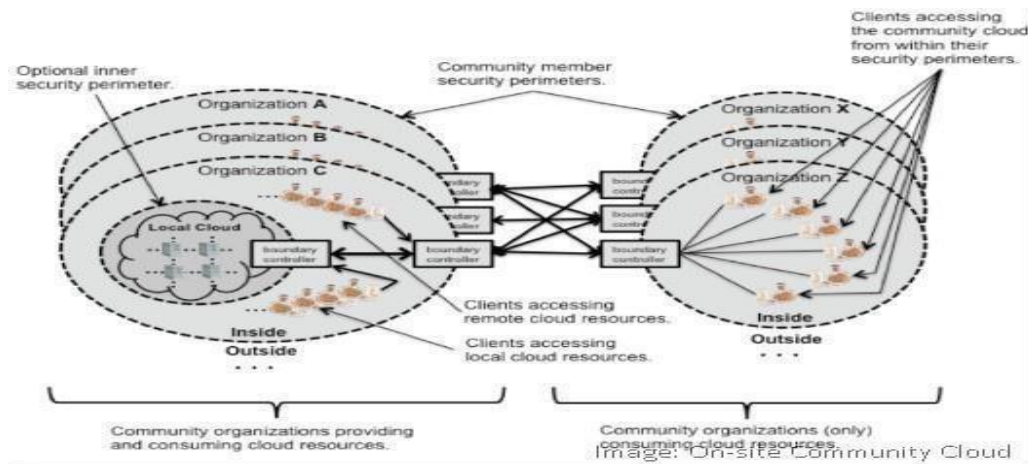- 



**Figure 1.2. Schematic Sketch of Public Cloud**

### 1.3.3  Public Cloud

This is the most common and widely recognized form of cloud computing. Here's a detailed overview:

- **Accessibility**: Available to the general public or a large industry group.
- **Location**: Typically off-site, hosted on the premises of the cloud service provider.
- **Ownership and Management**: Owned, operated, and managed by a third-party cloud service provider. Examples include Amazon Web Services (AWS), Microsoft Azure, and Google Cloud Platform (GCP).
- **Security and Control**: Security measures are implemented by the cloud service provider, but control over the infrastructure is limited for the user compared to private clouds. Users can manage their own data and applications, but the underlying infrastructure is controlled by the provider.

**Key Points**
- **Scalability**: Highly scalable due to the vast resources of cloud service providers. Users can easily scale their usage up or down based on demand.
- **Cost Efficiency**: Cost-effective as users pay only for the resources they use, eliminating the need for significant capital expenditure on hardware and maintenance.
- **Accessibility and Availability**: Services are accessible from anywhere with an internet connection, providing high availability and redundancy.

**Use Cases**
- **Startups and Small Businesses**: Can quickly deploy applications without the need for significant upfront investment in infrastructure.
- **Large Enterprises**: Use public clouds for non-sensitive workloads, development, and testing environments.
- **Software as a Service (SaaS) Providers**: Deliver applications to a broad audience over the internet.

**Examples of Public Cloud Services**

- **Infrastructure as a Service (IaaS)**: Provides virtualized computing resources over the internet. Examples: AWS EC2, Google Compute Engine.
- **Platform as a Service (PaaS)**: Offers hardware and software tools over the internet. Examples: Microsoft Azure, AWS Elastic Beanstalk.
- **Software as a Service (SaaS)**: Delivers software applications over the internet. Examples: Google Workspace, Salesforce.

Public clouds are the backbone of cloud computing, offering flexibility, scalability, and cost-efficiency to a wide range of users and organizations.

### 1.3.4  Hybrid Cloud

This combines multiple types of cloud deployment models. Here's a detailed overview:
- **Accessibility**: Involves a mix of private, public, and/or community clouds.
- **Location**: Can span on-site (private cloud) and off-site (public or community cloud) environments.
- **Ownership and Management**: Managed by the organization itself and/or third-party providers, depending on the specific setup of the hybrid cloud.
- **Security and Control**: Offers a balance of security and control by allowing organizations to keep sensitive data on private clouds while leveraging the scalability and cost-efficiency of public clouds for less sensitive workloads.

**Key Points**
- **Integration**: Hybrid clouds are bound together by standardized or proprietary technology that facilitates data and application portability. This enables seamless interaction and operation between different cloud environments.
- **Flexibility**: Allows organizations to optimize their resources by dynamically shifting workloads between private and public clouds based on demand, cost, and other factors.
- **Scalability and Efficiency**: Organizations can scale out to public clouds during peak demand (cloud bursting) and maintain critical operations on private clouds.

**Use Cases**
- **Disaster Recovery**: Use public clouds for backup and disaster recovery while keeping primary operations on a private cloud.
- **Load Balancing**: Utilize public clouds to handle traffic spikes and load balancing, ensuring high availability and performance.
- **Development and Testing**: Conduct development and testing in public clouds and deploy production environments on private clouds for enhanced security.

**Examples of Hybrid Cloud Scenarios**
- **Cloud Bursting**: An organization runs its baseline workload in a private cloud but "bursts" into a public cloud when demand spikes, ensuring resources are available without over-provisioning.
- **Data Processing**: Sensitive data is stored and processed on a private cloud, while large-scale data processing tasks are handled on a public cloud to take advantage of its computational power.

Hybrid clouds provide a versatile solution, combining the benefits of multiple cloud models to meet diverse business needs. They offer an ideal mix of security, scalability, and cost-efficiency, making them suitable for a wide range of applications and industries.
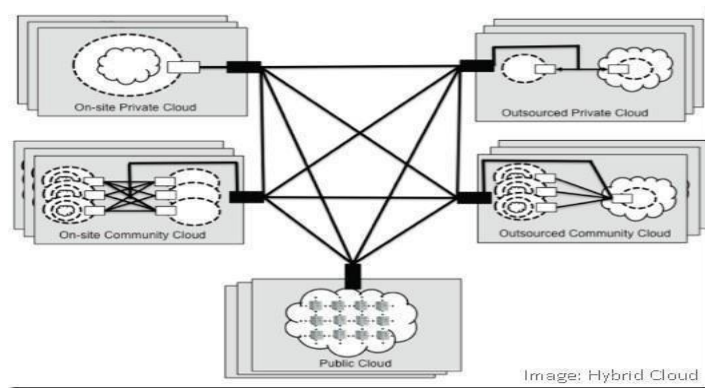
**Figure 1.3. Schematic Sketch of Hybrid Cloud**



**Figure 1.4. Cloud Computing Life Cycle**

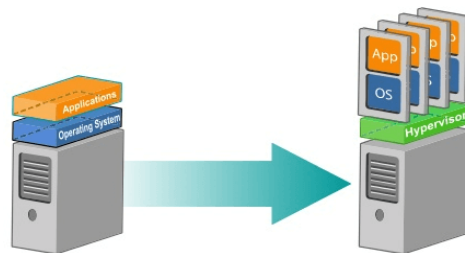## 1.4 Virtualization in cloud computing



**Figure 1.5. Concept of Virtualization**

Virtualization is a foundational technology for cloud computing. It enables the creation of virtual instances of computing resources, allowing multiple virtual machines (VMs) to run on a single physical machine. This abstraction of hardware resources provides the flexibility, scalability, and efficiency that are essential for cloud environments. Here's a detailed look at virtualization in cloud computing:

### 1.4.1 Key Concepts of Virtualization

1. **Hypervisor (Virtual Machine Monitor)**

- o **Type 1 (Bare-Metal Hypervisor)**: Runs directly on the physical hardware, providing better performance and efficiency. Examples include VMware ESXi, Microsoft Hyper-V, and Xen.
- o **Type 2 (Hosted Hypervisor)**: Runs on a host operating system, which then runs on the physical hardware. Examples include VMware Workstation and Oracle VirtualBox.

2. **Virtual Machines (VMs)**
   - o **Definition**: A VM is an emulation of a physical computer, running an operating system and applications.
   - o **Components**: VMs have their own virtual CPU, memory, storage, and network interfaces.
3. **Containers**
   - o **Definition**: Containers are lightweight, portable, and self-sufficient units that include an application and all its dependencies but share the host operating system's kernel.
   - o **Examples**: Docker, Kubernetes.
4. **Storage Virtualization**
   - o **Definition**: Abstracts physical storage resources to create a pool of storage that can be managed and allocated as needed.
   - o **Examples**: Storage Area Networks (SANs), Network Attached Storage (NAS).
5. **Network Virtualization**
   - o **Definition**: Creates a virtual network layer that abstracts physical network resources, enabling the creation of virtual networks.
   - o **Examples**: Virtual LAN (VLAN), Software-Defined Networking (SDN).
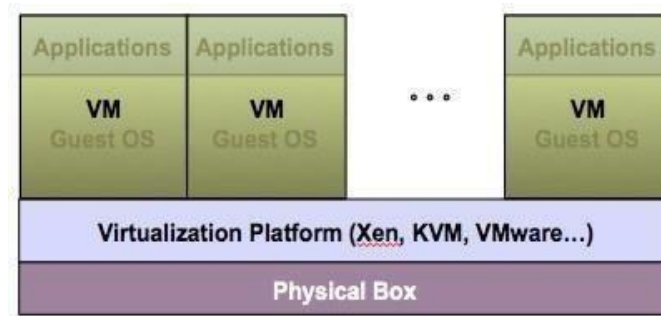


**Figure 1.6. Concept of Virtualization**
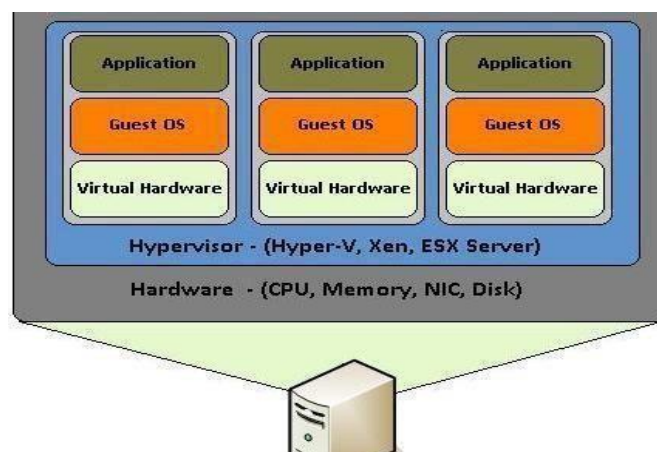
**1.5 Hypervisor**

**Figure 1.7. Schematic Sketch of Hypervisor**

The term hypervisor was first coined in 1956 by IBM
Hypervisor acts as a link between the hardware and the virtual environment and distributes the hardware resources such as CPU usage, memory allotment between the different virtual environments.
A hypervisor or virtual machine monitor (VMM) is computer software, firmware or hardware that creates and runs virtual machines. A computer on which a hypervisor runs one or more virtual machines is called a host machine, and each virtual machine is called a guest machine. A hypervisor is a hardware virtualization technique that allows multiple guest operating systems (OS) to run on a single host system at the same time. The guest OS shares the hardware of the host computer, such that each OS appears to have its own processor, memory and other hardware resources. A hypervisor is also known as a virtual machine manager (VMM).

### 1.5.1 Types of virtualization

Virtualization can be categorized into several types, each serving different purposes and providing various benefits in cloud computing and IT environments. Here are the main types of virtualization:

**1. Hardware Virtualization**
- **Definition**: Abstracts physical hardware resources to create virtual machines (VMs) that operate like physical computers.
- **Hypervisors**: The software layer that enables hardware virtualization.
  - **Type 1 (Bare-Metal Hypervisors)**: Runs directly on the physical hardware (e.g., VMware ESXi, Microsoft Hyper-V, Xen).
  - **Type 2 (Hosted Hypervisors)**: Runs on a host operating system (e.g., VMware Workstation, Oracle VirtualBox).
- **Use Cases**: Server consolidation, development and testing environments, running multiple operating systems on a single physical machine.

**2. Operating System Virtualization (Containerization)**
- **Definition**: Abstracts the operating system layer to run multiple isolated user-space instances (containers) on a single OS kernel.
- **Tools and Platforms**: Docker, Kubernetes, OpenVZ, LXC/LXD.
- **Use Cases**: Microservices architecture, rapid deployment and scaling of applications, isolation of applications without the overhead of full VMs.

**3. Network Virtualization**
- **Definition**: Abstracts physical network resources to create virtual networks that can be managed and optimized independently.
- **Components**: Virtual LANs (VLANs), Virtual Private Networks (VPNs), Software-Defined Networking (SDN).
- **Use Cases**: Network segmentation, enhanced security, flexible and programmable network management.

**4. Storage Virtualization**
- **Definition**: Abstracts physical storage resources to create a pool of storage that can be managed and allocated as needed.
- **Types**:
  - **Block-Level Storage Virtualization**: Abstracts blocks of storage, commonly used in SAN environments.
  - **File-Level Storage Virtualization**: Abstracts files and directories, commonly used in NAS environments.
- **Use Cases**: Simplified storage management, increased storage utilization, improved scalability and availability.

**5. Desktop Virtualization**
- **Definition**: Separates the desktop environment and associated applications from the physical client device.

- **Types**:
  - o **Virtual Desktop Infrastructure (VDI)**: Hosts desktop environments on VMs that run on centralized servers (e.g., VMware Horizon, Citrix Virtual Apps and Desktops).
  - o **Remote Desktop Services (RDS)**: Provides access to a desktop environment on a remote server.
- **Use Cases**: Centralized management of desktops, remote access, improved security and compliance.

## 6. Application Virtualization
- **Definition**: Abstracts applications from the underlying operating system, allowing them to run in isolated environments.
- **Tools and Platforms**: VMware ThinApp, Microsoft App-V, Citrix XenApp.
- **Use Cases**: Simplified application deployment and updates, application isolation, compatibility with different operating systems.

## 7. Data Virtualization
- **Definition**: Abstracts data from different sources to provide a unified view without the need for data replication or movement.
- **Tools and Platforms**: Denodo, IBM Data Virtualization, Red Hat JBoss Data Virtualization.
- **Use Cases**: Real-time data integration, business intelligence, data analysis from multiple sources.

## 8. Memory Virtualization
- **Definition**: Abstracts physical memory resources to create a pool of memory that can be dynamically allocated to applications as needed.
- **Techniques**: Virtual memory, memory paging, and swapping.
- **Use Cases**: Enhanced application performance, efficient memory utilization, support for large-scale applications.

## 1.6        Parallelization in cloud computing

Parallelization in cloud computing involves the simultaneous execution of multiple tasks to improve computational efficiency, reduce processing time, and enhance scalability. It is a critical concept for leveraging the full potential of cloud resources, particularly in applications that require significant computational power or handle large datasets. Here's a detailed look at parallelization in cloud computing:

**Key Concepts of Parallelization**

1. **Task Parallelism**
   - o **Definition**: Different tasks or processes run concurrently on multiple processors or cores.
   - o **Use Cases**: Independent tasks in a workflow, multi-threaded applications, distributed simulations.
2. **Data Parallelism**
   - o **Definition**: The same task is performed on different pieces of distributed data simultaneously.
   - o **Use Cases**: Large-scale data processing, machine learning training, scientific computing.
3. **Pipeline Parallelism**
   - o **Definition**: Different stages of a process are executed in parallel, similar to an assembly line.
   - o **Use Cases**: Streaming data processing, video processing, ETL (Extract, Transform, Load) workflows.

**Techniques and Tools for Parallelization**

1. **MapReduce**
   - **Definition**: A programming model for processing large datasets in parallel across a distributed cluster.
   - **Components**:
     - **Map**: Processes input data and produces intermediate key-value pairs.
     - **Reduce**: Aggregates intermediate results to produce final output.
   - **Platforms**: Apache Hadoop, Google Cloud Dataflow, Amazon EMR.
2. **Distributed Computing Frameworks**
   - **Apache Spark**: Provides in-memory data processing and is designed for large-scale data processing.
   - **Dask**: Parallel computing with Python, designed to parallelize NumPy, pandas, and scikit-learn operations.
   - **Flink**: Stream processing framework that supports batch processing as well.
3. **Cluster Management and Orchestration**
   - **Kubernetes**: Manages containerized applications across a cluster of machines, ensuring efficient resource utilization.
   - **Apache Mesos**: Manages resources and schedules tasks across a cluster.
4. **Grid Computing**
   - **Definition**: Uses a distributed network of loosely coupled computers to perform large-scale tasks.
   - **Use Cases**: Scientific research, complex simulations, large-scale computations.
5. **Serverless Computing**
   - **Definition**: Automatically scales resources in response to demand, enabling parallel execution of functions without managing servers.
   - **Platforms**: AWS Lambda, Azure Functions, Google Cloud Functions.

**Benefits of Parallelization in Cloud Computing**

1. **Scalability**
   - Efficiently scales applications by distributing tasks across multiple nodes or servers.
   - Handles large-scale data processing and high-performance computing workloads.
2. **Performance**
   - Reduces processing time by executing multiple tasks simultaneously.
   - Enhances throughput and responsiveness of applications.
3. **Cost Efficiency**
   - Optimizes resource utilization, reducing the need for over-provisioning.
   - Pay-as-you-go pricing models in the cloud allow cost savings by leveraging parallel processing.
4. **Flexibility**
   - Supports a wide range of applications, from data analytics to machine learning and real-time processing.
   - Enables dynamic allocation of resources based on workload demands.

**Challenges of Parallelization**

1. **Synchronization and Coordination**
   - Ensuring tasks are properly synchronized to avoid data inconsistencies and deadlocks.
   - Managing dependencies between tasks.

2. **Resource Management**
    o Efficiently allocating resources to avoid bottlenecks.
    o Balancing load across distributed systems.
3. **Fault Tolerance**
    o Handling failures in a distributed environment to ensure reliability and availability.
    o Implementing recovery mechanisms.
4. **Complexity**
    o Developing and debugging parallel applications can be more complex than serial ones.
    o Requires expertise in parallel programming and distributed systems.


**Cloud Resource Management**

Cloud resource management involves efficiently allocating, monitoring, and optimizing cloud resources to ensure cost-effectiveness, performance, and reliability. Effective resource management is essential for maximizing the benefits of cloud computing, including scalability, flexibility, and cost savings. Here's an in-depth look at cloud resource management:

**Key Aspects of Cloud Resource Management**
1. **Resource Provisioning**
    o **Automated Provisioning**: Automatically allocating resources based on predefined policies or real-time demand.
    o **Manual Provisioning**: Manually allocating resources based on user requests or administrative decisions.
    o **Elasticity**: The ability to scale resources up or down based on workload requirements.
2. **Resource Monitoring**
    o **Performance Metrics**: Tracking CPU, memory, storage, and network utilization.
    o **Health Monitoring**: Monitoring the status of resources to detect failures or performance issues.
    o **Usage Analytics**: Analyzing resource usage patterns to identify inefficiencies and optimize utilization.
3. **Resource Optimization**
    o **Load Balancing**: Distributing workloads across multiple resources to ensure optimal performance and avoid bottlenecks.
    o **Cost Optimization**: Identifying and eliminating unnecessary expenses by rightsizing resources and utilizing cost-effective options.
    o **Capacity Planning**: Forecasting future resource needs based on historical data and anticipated growth.
4. **Resource Governance**
    o **Policy Management**: Implementing policies to control resource allocation, usage, and access.
    o **Compliance**: Ensuring resource usage adheres to regulatory requirements and organizational standards.
    o **Access Control**: Managing permissions and access rights to resources to ensure security and compliance.

**Tools and Technologies for Cloud Resource Management**
1. **Cloud Management Platforms (CMPs)**
    o **Features**: Provide a unified interface for managing resources across multiple cloud providers.
    o **Examples**: VMware vRealize, Cisco CloudCenter, Scalr.
2. **Native Cloud Tools**
    o **AWS**: AWS CloudFormation, AWS CloudWatch, AWS Auto Scaling, AWS

Trusted Advisor.
- o **Azure**: Azure Resource Manager, Azure Monitor, Azure Advisor, Azure Cost Management.
- o **Google Cloud**: Google Cloud Deployment Manager, Stackdriver, Google Kubernetes Engine (GKE) Autoscaler, Google Cloud Cost Management.
3. **Third-Party Monitoring and Optimization Tools**
   - o **Performance Monitoring**: Datadog, New Relic, AppDynamics.
   - o **Cost Management**: CloudHealth, Cloudability, Spot by NetApp.
   - o **Security and Compliance**: Dome9, CloudCheckr, Palo Alto Networks Prisma Cloud.

**Strategies for Effective Cloud Resource Management**
1. **Automation**
   - o **Infrastructure as Code (IaC)**: Using code to manage and provision cloud resources (e.g., Terraform, AWS CloudFormation).
   - o **Automated Scaling**: Implementing auto-scaling policies to dynamically adjust resources based on demand.
2. **Monitoring and Analytics**
   - o **Real-Time Monitoring**: Continuously monitoring resource performance and health.
   - o **Predictive Analytics**: Using machine learning and historical data to predict future resource needs and optimize capacity planning.
3. **Cost Management**
   - o **Rightsizing**: Continuously evaluating and adjusting resource sizes to match current workload requirements.
   - o **Reserved Instances and Savings Plans**: Leveraging long-term commitment options for cost savings.
   - o **Tagging and Resource Grouping**: Implementing a tagging strategy to track and manage costs by project, department, or application.
4. **Governance and Compliance**
   - o **Policy Enforcement**: Using policies to enforce best practices and compliance requirements.
   - o **Access Management**: Implementing role-based access control (RBAC) and least privilege principles to secure resources.
5. **Continuous Improvement**
   - o **Feedback Loops**: Regularly reviewing performance and cost metrics to identify areas for improvement.
   - o **Iterative Optimization**: Continuously refining resource management practices based on feedback and evolving requirements.

## 1.7 Optimal allocation of cloud models

Optimal allocation of cloud models involves choosing and distributing cloud resources and services in a manner that maximizes efficiency, performance, and cost-effectiveness while meeting the specific needs of the organization. This involves deciding which cloud deployment models (public, private, hybrid, or multi-cloud) and service models (IaaS, PaaS, SaaS) to use for various applications and workloads. Here's a comprehensive look at the optimal allocation of cloud models:

**Factors Influencing Cloud Model Allocation**

1. **Workload Characteristics**
   - o **Performance Requirements**: High-performance applications may benefit from private clouds or dedicated resources.
   - o **Scalability Needs**: Applications with variable demand can leverage the scalability of public clouds.

- o **Data Sensitivity**: Sensitive data might require private clouds or specific compliance configurations.
2. **Cost Considerations**
   - o **Budget Constraints**: Public clouds often provide cost advantages through pay-as-you-go pricing.
   - o **Total Cost of Ownership (TCO)**: Includes not just direct costs but also indirect costs like maintenance and management.
3. **Compliance and Security**
   - o **Regulatory Requirements**: Some industries require data to be stored in specific geographic locations or under strict compliance guidelines.
   - o **Security Policies**: Organizations may prefer private or hybrid clouds to maintain greater control over security.
4. **Existing Infrastructure**
   - o **Legacy Systems**: Integrating with existing on-premises infrastructure may necessitate a hybrid approach.
   - o **Interoperability**: Ensuring that new cloud services can seamlessly integrate with existing systems.
5. **Operational Flexibility**
   - o **Deployment Speed**: Public clouds can quickly provision resources for rapid deployment.
   - o **Customizability**: Private clouds offer greater customization options for specific organizational needs.
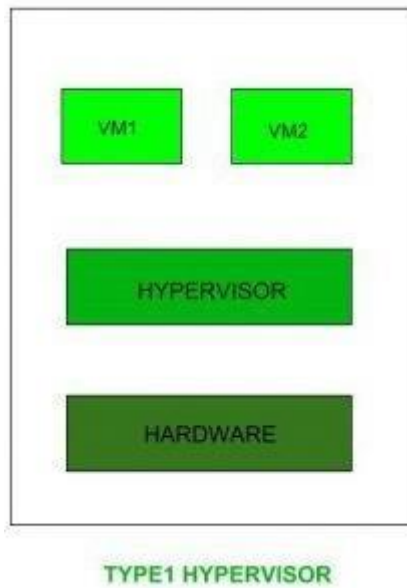
**Hypervisor**

A hypervisor, also known as a virtual machine monitor (VMM), is a crucial component in virtualization technology that enables multiple operating systems (OS) to run concurrently on a single physical machine. It abstracts the underlying hardware resources and creates virtual environments, known as virtual machines (VMs), where each VM operates as if it were a standalone computer with its own CPU, memory, storage, and network interfaces. Here's a detailed overview of hypervisors and their role in virtualization:

## Types of Hypervisors

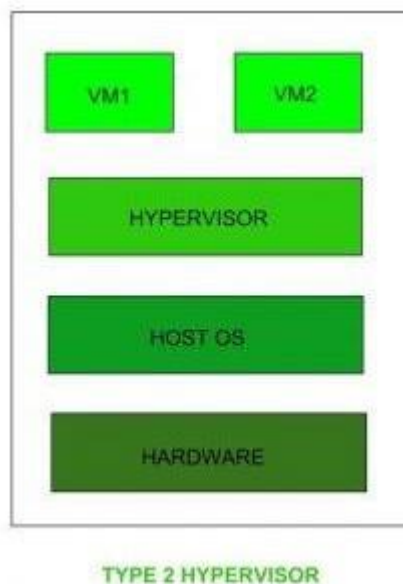Hypervisors are classified into two main types based on their architecture and deployment:

1. **Type 1 Hypervisor (Bare-Metal Hypervisor)**

- **Architecture**: Runs directly on the physical hardware (bare-metal).
- **Example**: VMware ESXi, Microsoft Hyper-V, Citrix XenServer, KVM (Kernel-based Virtual Machine).
- **Advantages**:
  - o Efficient performance as it directly interacts with hardware resources.
  - o Minimal overhead compared to Type 2 hypervisors.
  - o Typically used in enterprise data centers for server virtualization.

TYPE1 HYPERVISOR

2. **Type 2 Hypervisor (Hosted Hypervisor)**

- **Architecture**: Runs on top of a host operating system.
- **Example**: VMware Workstation, Oracle VirtualBox, Parallels Desktop.
- **Advantages**:
    - Easy installation and setup, suitable for desktop virtualization and testing environments.
    - Can run on a variety of host operating systems (Windows, macOS, Linux).



TYPE 2 HYPERVISOR

## Functions and Features

- **Resource Virtualization**: Hypervisors abstract physical CPU, memory, storage, and network resources into virtual resources allocated to VMs.
- **Isolation**: Provides strong isolation between VMs, ensuring that each VM operates independently without affecting others.

- **Resource Management**: Manages and allocates resources dynamically based on VM workload demands.
- **Virtual Machine Lifecycle Management**: Handles VM creation, provisioning, migration, and deletion.
- **Security**: Enforces access controls and security policies to protect VMs and their data.
- **Live Migration**: Allows VMs to be migrated between physical hosts with minimal downtime, facilitating load balancing and maintenance activities.
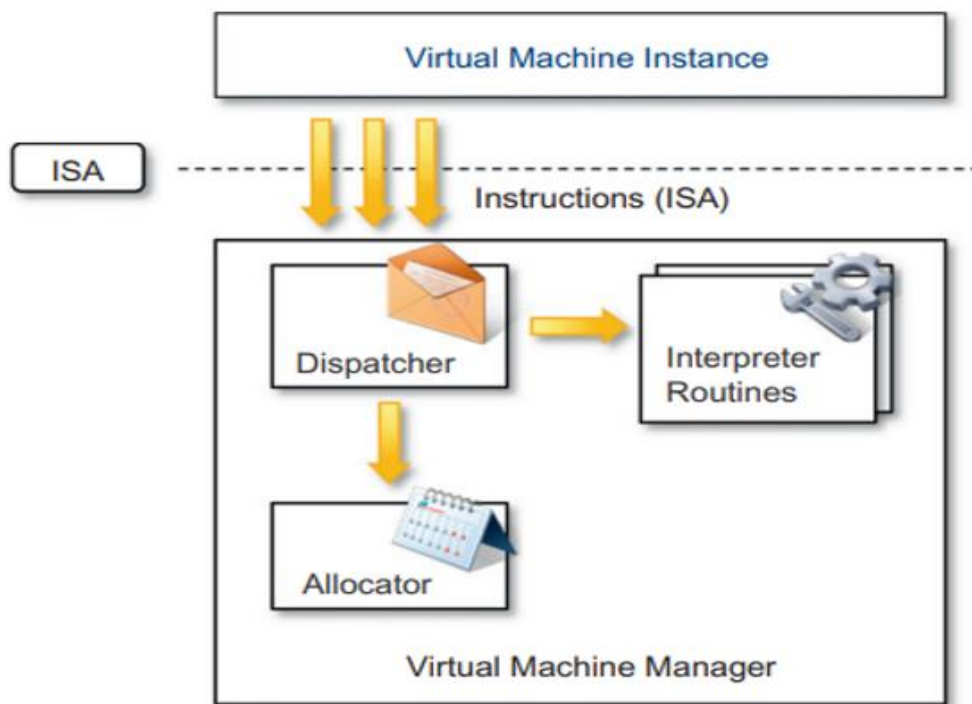
**HYPERVISOR REFERENCE MODEL**



**Fig. 1.8 Hypervisor reference architecture**

There are 3 main modues coordinate in order to emulate the undrelying hardware:
1. Dispatcher
2. Allocator
3. Interpreter

**DISPATCHER:**
The dispatcher behaves like the entry point of the monitor and reroutes the instructions of the virtualmachine instance to one of the other two modules.

**ALLOCATOR:**
The allocator is responsible for deciding the system resources to be provided to the virtual machine instance. It means whenever virtual machine tries to execute an instruction that results in changing the machine resources associated with the virtual machine, the allocator is invoked by the dispatcher.

**INTERPRETER:**
The interpreter module consists of interpreter routines. These are executed, whenever virtual machineexecutes a priviliged instruction.

**Virtualization Case Studies:**
**Virtualization Structures, Tools and Mechanisms**
In general, there are three typical classes of VM architecture. The virtualization layer is responsible for converting portions of the real hardware into virtual hardware. Therefore, different operating systems such as Linux and Windows can run on the same physical machine, simultaneously. Depending on the position of the virtualization layer, there are several classes of VM architectures, namely the hypervisor architecture, paravirtualization and host based virtualization. The hypervisor is also known as the VMM (Virtual Machine Monitor). They both perform the same virtualization operations.

## 1.9 Hypervisor and Xen architecture

The hypervisor supports hardware level virtualization on bare metal devices like CPU, memory, disk and network interfaces. The hypervisor software sits directly between the physical hardware and its OS. This virtualization layer is referred to as either the VMM or the hypervisor. The hypervisor provides hypercalls for the guest OSes and applications. Depending on the functionality, a hypervisor can assume micro kernel architecture like the Microsoft Hyper-V. It can assume monolithic hypervisor architecture like the VMware ESX for server virtualization. A micro kernel hypervisor includes only the basic and unchanging functions (such as physical memory management and processor scheduling). The device drivers and other changeable components are outside the hypervisor. A monolithic hypervisor implements all the aforementioned functions, including those of the device drivers. Therefore, the size of the hypervisor code of a micro-kernel hypervisor is smaller than that of a monolithic hypervisor. Essentially, a hypervisor must be able to convert physical devices into virtual resources dedicated for the deployed VM to use.

## 1.9.1 Xen architecture

Xen is an open source hypervisor program developed by Cambridge University. Xen is a microkernel hypervisor, which separates the policy from the mechanism. The Xen hypervisor implements all the mechanisms, leaving the policy to be handled by Domain 0. Figure 2.9 shows architecture of Xen hypervisor. Xen does not include any device drivers natively. It just provides a mechanism by which a guest OS can have direct access to the physical devices. As a result, the size of the Xen hypervisor is kept rather small. Xen provides a virtual environment located between the hardware and the OS.
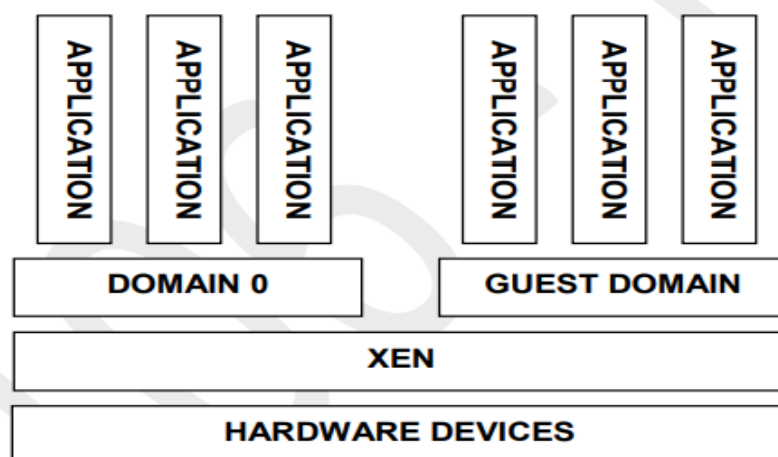


**Fig. 1.9 Xen domain 0 for control and I/O & guest domain for user applications.**

The core components of a Xen system are:
- ✓ the hypervisor,
- ✓ kernel, and

✓ Applications.

The organization of the three components is important. Like other virtualization systems, many guest OSes can run on top of the hypervisor. However, not all guest OSes are created equal, and one in particular controls the others.

The guest OS, which has control ability, is called Domain 0, and the others are called Domain U.  Domain 0 is a privileged guest OS of Xen. It is first loaded when Xen boots without any file system drivers being available. Domain 0 is designed to access hardware directly and manage devices. Therefore, one of the responsibilities of Domain 0 is to allocate and map hardware resources for the guest domains (the Domain U domains). For example, Xen is based on Linux and its security level is C2. Its management VM is named Domain 0 which has the privilege to manage other VMs implemented on the same host. If Domain 0 is compromised, the hacker can control the entire system. So, in the VM system, security policies are needed to improve the security of Domain 0. Domain 0, behaving as a VMM, allows users to create, copy, save, read, modify, share, migrate and roll back VMs as easily as manipulating a file, which flexibly provides tremendous benefits for users.

Figure describes the architecture of Xen and its mapping onto a classic x86 privilege model. A Xen-based system is managed by the Xen hypervisor, which runs in the highest privileged mode and controls the access of guest operating system to the underlying hardware.
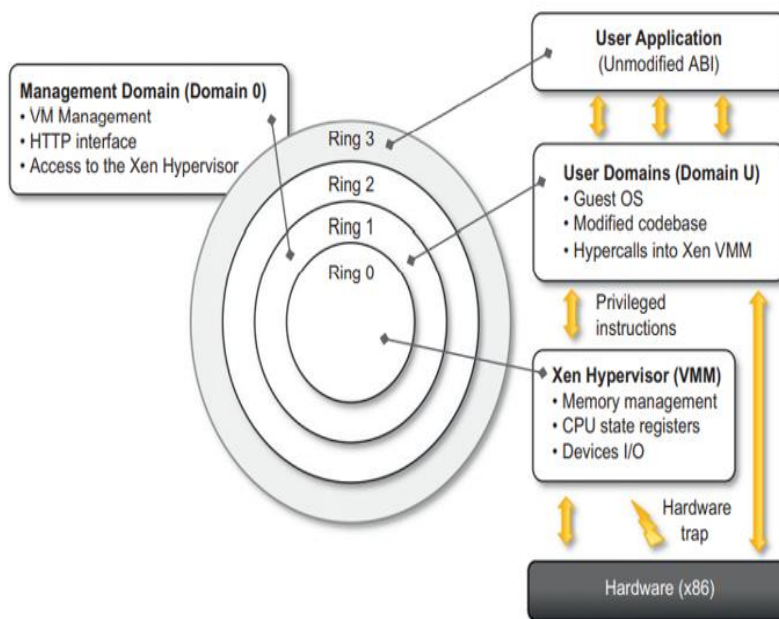


**Fig. 1.10 Xen architecture and OS management**

### 1.9 VMware: full virtualization
VMware's technology is based on the concept of full virtualization, where the underlying hardware is replicated and made available to the guest operating system, which runs unaware of such abstraction layers and does not need to be modified. VMware implements full virtualization either in the desktop environment, by means of Type II hypervisors, or in the server environment, by means of Type I hypervisors.

### 1.10 Microsoft Hyper-V
Hyper-V is an infrastructure virtualization solution developed by Microsoft for server

virtualization.

It uses a hypervisor-based approach to hardware virtualization, which leverages several techniques to support a variety of guest operating systems. Hyper-V is currently shipped as a component of Windows Server 2008 R2 that installs the hypervisor as a role within the server.

### 1.10.1 Architecture

Hyper-V supports multiple and concurrent execution of guest operating systems by means of partitions. A partition is a completely isolated environment in which an operating system is installed and run. Figure below provides an overview of the architecture of Hyper-V. Despite its straightforward installation as a component of the host operating system, Hyper-V takes control of the hardware, and the host operating system becomes a virtual machine instance with special privileges, called the parent partition.
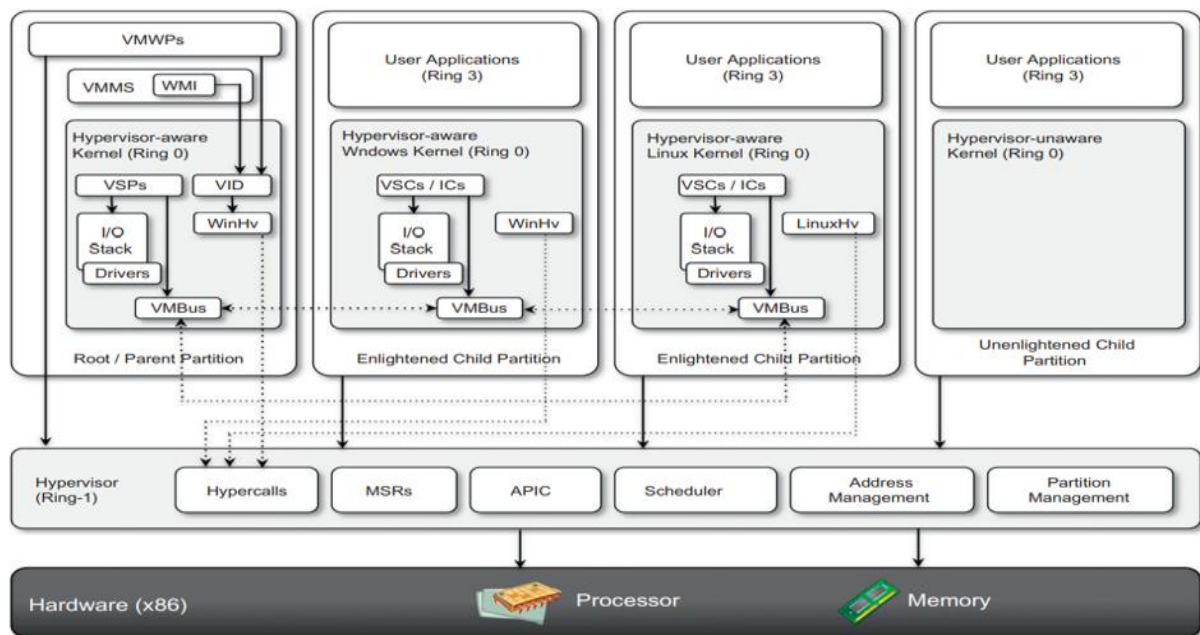


**Fig. 1.11 Hyper V architecture**

The parent partition (also called the root partition) is the only one that has direct access to the hardware. It runs the virtualization stack, hosts all the drivers required to configure guest operating systems, and creates child partitions through the hypervisor. Child partitions are used to host guest operating systems and do not have access to the underlying hardware, but their interaction with it is controlled by either the parent partition or the hypervisor itself.

**Hypervisor** The hypervisor is the component that directly manages the underlying hardware (processors and memory). It is logically defined by the following components:

- **Hypercalls interface**
- Memory service routines (MSRs)
- Advanced programmable interrupt controller (APIC)
- Scheduler
- Address manager
- Partition manager

**Hypercalls interface**. This is the entry point for all the partitions for the execution of sensitive instructions. This is an implementation of the paravirtualization approach already discussed with Xen. This interface is used by drivers in the partitioned operating system to contact the hypervisor using the standard Windows calling convention.

**Memory service routines (MSRs)**. These are the set of functionalities that control the memory and its access from partitions. By leveraging hardware-assisted virtualization, the hypervisor uses the Input/Output Memory Management Unit (I/O MMU or IOMMU) to fast-track access to devices from partitions by translating virtual memory addresses.

**Advanced programmable interrupt controller (APIC)**. This component represents the interrupt controller, which manages the signals coming from the underlying hardware when some event occurs (timer expired, I/O ready, exceptions and traps). Each virtual processor is equipped with a synthetic interrupt controller (SynIC), which constitutes an extension of the local APIC. The hypervisor is responsible of dispatching, when appropriate, the physical interrupts to the synthetic interrupt controllers.

**Scheduler**. This component schedules the virtual processors to run on available physical processors. The scheduling is controlled by policies that are set by the parent partition.

**Address manager**. This component is used to manage the virtual network addresses that are allocated to each guest operating system.

**Partition manager.** This component is in charge of performing partition creation, finalization, destruction, enumeration, and configurations.

The hypervisor runs in Ring -1 and therefore requires corresponding hardware technology that enables such a condition. By executing in this highly privileged mode, the hypervisor can support legacy operating systems that have been designed for x86 hardware. Operating systems of newer generations can take advantage of the new specific architecture of Hyper-V especially for the I/O operations performed by child partitions.