

Biyani's Think Tank

Concept Based Notes

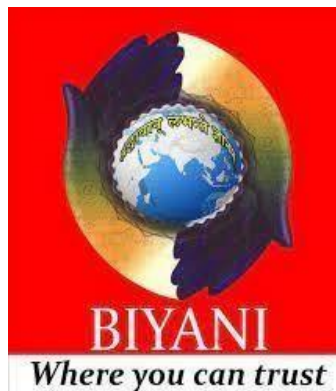
Cloud Computing

Mr. Himanshu Mathur

Asst. Professor (Dept. of IT)

Biyani Girls College, Jaipur

BCA VI Semester



Published by :

Think Tanks

Biyani Group of Colleges

Concept & Copyright :

Biyani Shikshan Samiti Sector-3, Vidhyadhar Nagar, Jaipur-302 023 (Rajasthan)

Ph : 0141-2338371, 2338591-95 Fax : 0141-2338007

E-mail : acad@biyanicolleges.org

Website : www.gurukpo.com; www.biyanicolleges.org

ISBN :

Edition: 2025

While every effort is taken to avoid errors or omissions in this Publication, any mistake or omission that may have crept in is not intentional. It may be taken note of that neither the publisher nor the author will be responsible for any damage or loss of any kind arising to

Leaser Type Setted by :

Biyani College Printing Department

Preface

I am glad to present this book, especially designed to serve the needs of the students. The book has been written keeping in mind the general weakness in understanding the fundamental concepts of the topics. The book is self-explanatory and adopts the —Teach Yourself style. It is based on question- answer pattern. The language of book is quite easy and understandable based on scientific approach.

Any further improvement in the contents of the book by making corrections, omission and inclusion is keen to be achieved based on suggestions from the readers for which the author shall be obliged.

I acknowledge special thanks to Mr. Rajeev Biyani, *Chairman* & Dr. Sanjay Biyani, *Director (Acad.)* Biyani Group of Colleges, who are the backbones and main concept provider and also have been constant source of motivation throughout this Endeavour. They played an active role in coordinating the various stages of this Endeavour and spearheaded the publishing work.

I look forward to receiving valuable suggestions from professors of various educational institutions, other faculty members and students for improvement of the quality of the book. The reader may feel free to send in their comments and suggestions to the under mentioned address.

Author

Unit-I

Introduction of Cloud Computing: Definition, Historical Developments, Enabling Technology, Vision, Essential Characteristics of Cloud Computing , Components of Cloud Computing. Challenges and Approaches of Migration into Cloud, Cloud Applications:– Health care, CRM and ERP, Social Networking, Media Applications and Multiplayer Online Gaming. Benefits: For the Market, Enterprise, End user and Individuals.

Unit-II

Cloud Computing Architecture : Introduction, Cloud Reference Model, Architecture, Infrastructure / Hardware as a Service, Platform as a Service, Software as a Service, Types of Clouds, Public Clouds, Private Clouds, Hybrid Clouds, Community Clouds, Economics of the Cloud, Open Challenges, Cloud Interoperability and Standards, Scalability and Fault Tolerance. Parallel and distributed Computing-MapReduce, High level Language for Cloud, Service Oriented Computing.

Unit-III

Virtualization: Introduction, Characteristics of Virtualized Environment, Taxonomy of Virtualization Techniques, Virtualization and Cloud computing, Virtualization: of CPU, Memory, I/O Devices, Server , Desktop, Network, and data-center. Pros and Cons of Virtualization, Technology Examples-VMware and Microsoft Hyper-V,KVM, Xen. Introduction of Cloud security services, Design Principles, Policy Implementation, Cloud Computing Security Challenges, Cloud Computing Security Architecture. Cloud Security technologies to secure the data in Private and Public. Security Concerns. Risk Mitigation, Understanding and Identification of Threats in Cloud, SLA-Service Level Agreements.

Unit-IV

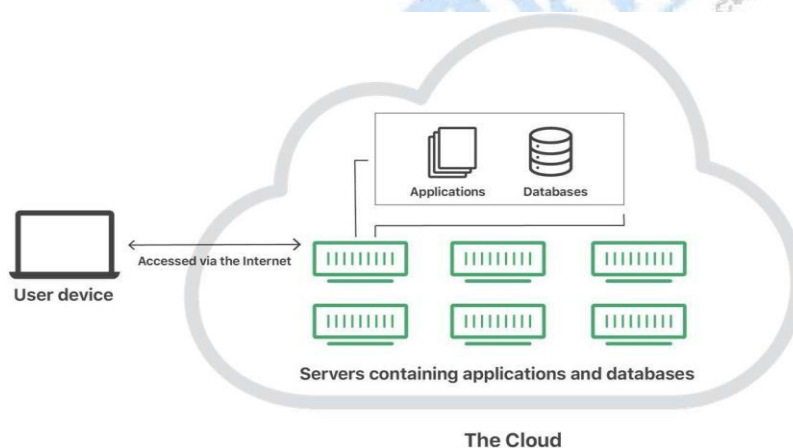
Cloud Platforms in Industry: Amazon Web Services- Compute Services, Storage Services, Communication Services and Additional Services. Google AppEngine-Architecture and Core Concepts, Application Life-Cycle, cost model.. Recommended

UNIT – I

Cloud Computing Fundamentals: Definition of Cloud computing, Roots of Cloud Computing, Layers and Types of Clouds, Desired Features of a Cloud, Cloud Infrastructure Management, Infrastructure as a Service Providers, Platform as a Service Providers. Computing Paradigms: High-Performance Computing, Parallel Computing, Distributed Computing, Cluster Computing, Grid Computing.

Introduction to Cloud Computing:

What is the cloud: "The cloud" refers to servers that are accessed over the Internet, and the software and databases that run on those servers. Cloud servers are in data centres all over the world. By using cloud computing, users and companies do not have to manage physical servers themselves or run software applications on their own machines.



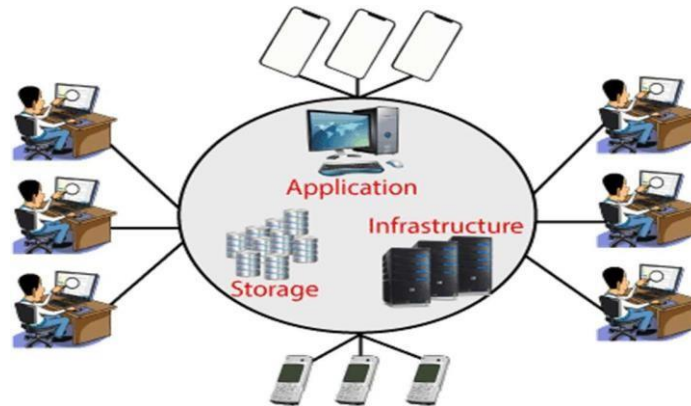
The cloud enables users to access the same files and applications from almost any device, because the computing and storage takes place on servers in a data centre, instead of locally on the user device. Therefore, a user can log into their Instagram account on a new phone after their old phone breaks and still find their old account in place, with all their photos, videos, and conversation history. It works the same way with cloud email providers like Gmail or Microsoft Office 365, and with cloud storage providers like Dropbox or Google Drive.

For businesses, switching to cloud computing removes some IT costs and overhead: for instance, they no longer need to update and maintain their own servers, as the cloud vendor they are using will do that. This especially makes an impact for small businesses that may not have been able to afford their own internal infrastructure but can outsource their infrastructure needs affordably via the cloud. The cloud can also make it easier for companies to operate internationally, because employees and customers can access the same files and applications from any location.

Definition of Cloud Computing:

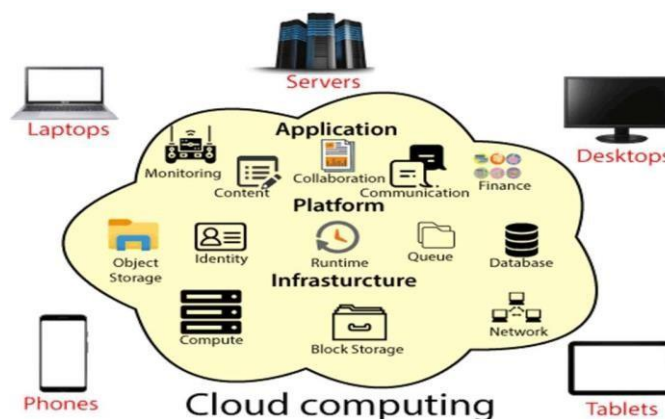
The term —Cloud Computing‖ refers to services provided by the cloud that is responsible for delivering of computing services such as servers, storage, databases, networking, software, analytics, intelligence, and more, over the Cloud (Internet).

Cloud computing applies a virtualized platform with elastic resources on demand by provisioning hardware, software, and data sets dynamically



Cloud Computing provides an alternative to the on-premises data center. With an on-premises data center, we must manage everything, such as purchasing and installing hardware, virtualization, installing the operating system, and any other required applications, setting up the network, configuring the firewall, and setting up storage for data. After doing all the set-up, we become responsible for maintaining it through its entire lifecycle.

However, if we choose Cloud Computing, a cloud vendor is responsible for the hardware purchase and maintenance. They also provide a wide variety of software and platform as a service. We can take any required services on rent. The cloud computing services are charged based on usage.



The cloud environment provides an easily accessible online portal that makes handy for the user to manage the compute, storage, network, and application resources. Some of the cloud service providers are in the following figure.



Advantages of cloud computing:

1. **Cost:** It reduces the huge capital costs of buying hardware and software.
2. **Speed:** Resources can be accessed in minutes, typically within a few clicks.
3. **Scalability:** We can increase or decrease the requirement of resources according to the business requirements.
4. **Productivity:** While using cloud computing, we put less operational effort. We do not need to apply patching, as well as no need to maintain hardware and software. So, in this way, the IT team can be more productive and focus on achieving business goals.
5. **Reliability:** Backup and recovery of data are less expensive and extremely fast for business continuity.
6. **Security:** Many cloud vendors offer a broad set of policies, technologies, and controls that strengthen our data security.

Cloud computing shares characteristics with:

1. **Client-server model**—*Client-server computing* refers broadly to any distributed application that distinguishes between service providers (servers) and service requestors (clients).
2. **Grid computing**—A form of distributed and parallel computing, whereby a 'super and virtual computer' is composed of a cluster of networked, loosely coupled computers acting in concert to perform very large tasks.
3. **Fog computing**—Distributed computing paradigm that provides data, compute, storage and application services closer to the client or near-user edge devices, such as network routers. Furthermore, fog computing handles data at the network level, on smart devices and on the end-user client-side (e.g., mobile devices), instead of sending data to a remote location for processing.
4. **Mainframe computer**—Powerful computers used mainly by large organizations for critical applications, typically bulk data processing such as census; industry and consumer statistics; police and secret intelligence services; enterprise resource planning; and financial transaction processing.
5. **Utility computing**—The packaging of computing resources, such as computation and storage, as a metered service similar to a traditional public utility, such as electricity.

6. **Peer-to-peer**—A distributed architecture without the need for central coordination. Participants are both suppliers and consumers of resources (in contrast to the traditional client-server model).
7. **Green computing**—Study and practice of environmentally sustainable computing or IT.
8. **Cloud sandbox**—A live, isolated computer environment in which a program, code or file can run without affecting the application in which it runs.

Characteristics of Cloud Computing

1. Agility for organizations
2. Cost reductions, Centralization of infrastructure in locations with lower costs.
3. Device and location independence, which means no maintenance, required.
4. Pay-per-use means utilization and efficiency improvements for systems that are often only 10–20% utilized.
5. Performances are being monitored by IT experts i.e., from the service provider end.
6. Productivity increases which results in multiple users who can work on the same data simultaneously.
7. Time may be saved as information does not need to be re-entered when fields are matched
8. Availability improves with the use of multiple redundant sites
9. Scalability and elasticity via dynamic ("on-demand") provisioning of resources on a fine-grained, self-service basis in near real-time without users having to engineer for peak loads.
10. Self-service interface.
11. Resources that are abstracted or virtualized.
12. Security can improve due to centralization of data

The National Institute of Standards and Technology's definition of cloud computing identifies "five essential characteristics":

1. On-demand self-service.
2. Broad network access.
3. Resource pooling.
4. Rapid elasticity.
5. Measured service.

ROOTS OF CLOUD COMPUTING

We can track the roots of clouds computing by observing the advancement of several technologies, especially in hardware (virtualization, multi-core chips), Internet technologies (Web services, service-oriented architectures, Web 2.0), distributed computing (clusters, grids), and systems management (autonomic computing, data center automation). Figure 1.1 shows the convergence of technology fields that

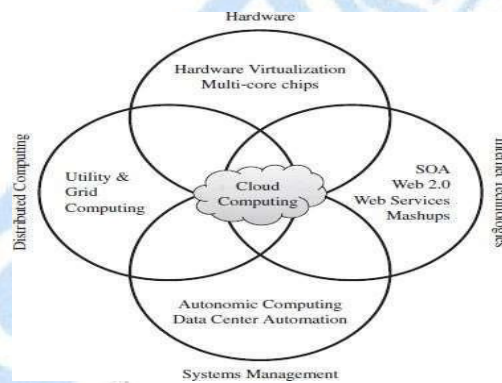
significantly advanced and contributed to the advent of cloud computing.

The emergence of cloud computing itself is closely linked to the maturity of such technologies. We present a closer look at the technologies that form the base of cloud computing, with the aim of providing a clearer picture of the cloud ecosystem.

1. From Mainframes to Clouds

We are currently experiencing a switch in the IT world, from in-house generated computing power into utility-supplied computing resources delivered over the Internet as Web services. This trend is like what occurred about a century ago when factories, which used to generate their own electric power, realized that it is was cheaper just plugging their machines into the newly formed electric power grid.

Computing delivered as a utility can be defined as —on demand delivery of infrastructure, applications, and business processes in a security-rich, shared, scalable, and based computer environment over the Internet for a feel



2. SOA, Web Services, Web 2.0, and Mashups

The emergence of Web services (WS) open standards has significantly contributed to advances in the domain of software integration. Web services can glue together applications running on different messaging product platforms, enabling information from one application to be made available to others, and enabling internal applications to be made available over the Internet.

WS standards have been created on top of existing ubiquitous technologies such as HTTP and XML, thus providing a common mechanism for delivering services, making them ideal for implementing a service-oriented architecture (SOA). The purpose of a SOA is to address requirements of loosely coupled, standards-based, and protocol-independent distributed computing. In a SOA, software resources are packaged as —services, which are well-defined, self-contained modules that provide standard business functionality and are independent of the state or context of other services. Services are described in a standard definition language and have a published interface

Many service providers, such as Amazon, delicious, Facebook, and Google, make their service APIs publicly accessible using standard protocols such as SOAP and REST. Consequently, one can put an idea of a fully functional Web application into practice just by gluing pieces with few lines of code.

1. Grid Computing

Grid computing enables aggregation of distributed resources and transparently access to them. Most production grids such as TeraGrid and EGEE seek to share compute and storage resources distributed across different administrative domains, with their focus being speeding up a broad range of scientific applications, such as climate modeling, drug design, and protein analysis.

A key aspect of the grid vision realization has been building standard Web services-based protocols that allow distributed resources to be —discovered, accessed, allocated, monitored, accounted for, and billed for, etc., and in general managed as a single virtual system.¶ The Open Grid Services Architecture (OGSA) addresses this need for standardization by defining a set of core capabilities and behaviors that address key concerns in grid systems.

Virtualization technology has been identified as the perfect fit to issues that have caused frustration when using grids, such as hosting many dissimilar software applications on a single physical platform. In this direction, some research projects (e.g., Globus Virtual Workspaces) aimed at evolving grids to support an additional layer to virtualize computation, storage, and network resources.

2. Utility Computing

With increasing popularity and usage, large grid installations have faced new problems, such as excessive spikes in demand for resources coupled with strategic and adversarial behavior by users. Initially, grid resource management techniques did not ensure fair and equitable access to resources in many systems. Traditional metrics (throughput, waiting time, and slowdown) failed to capture the more subtle requirements of users. There were no real incentives for users to be flexible about resource requirements or job deadlines, nor provisions to accommodate users with urgent work.

In utility computing environments, users assign a —utility¶ value to their jobs, where utility is a fixed or time-varying valuation that captures various QoS constraints (deadline, importance, satisfaction). The valuation is the amount they are willing to pay a service provider to satisfy their demands. The service providers then attempt to maximize their own utility, where said utility may directly correlate with their profit. Providers can choose to prioritize high yield (i.e., profit per unit of resource) user jobs, leading to a scenario where shared systems are viewed as a marketplace, where users compete for resources based on the perceived utility or value of their jobs. Further information and comparison of these utility computing environments are available in an extensive survey of these platforms

3. Hardware Virtualization

Cloud computing services are usually backed by large-scale data centers composed of thousands of computers. Such data centers are built to serve many users and host many disparate applications. For this purpose, hardware virtualization can be considered as a perfect fit to overcome most operational issues of data center building and maintenance.

The idea of virtualizing a computer system resources, including processors, memory, and I/O devices, has been well established for decades, aiming at improving sharing and utilization of computer systems. Hardware virtualization allows running multiple operating systems and software stacks on a single physical platform. As depicted in Figure 1.2, a software layer, the virtual machine monitor (VMM), also called a hypervisor, mediates access to the physical hardware presenting to each guest operating system a virtual machine (VM), which is a set of virtual platform interfaces

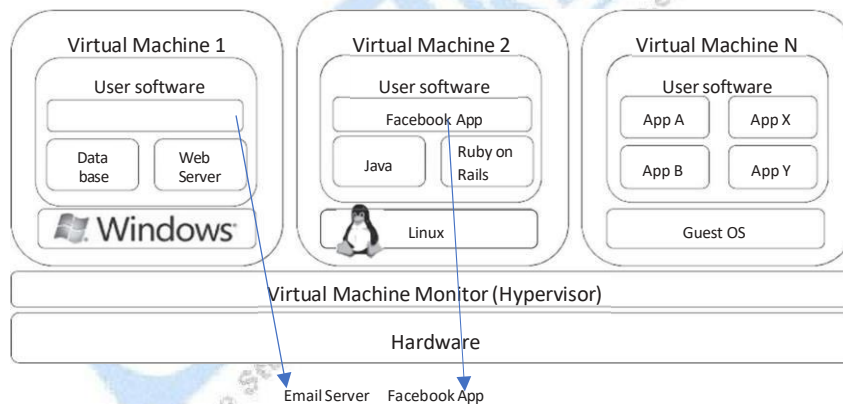


FIGURE 1.2. A hardware virtualized server hosting three virtual machines, each one running distinct operating system and user level software stack.

Several VMM platforms exist that are the basis of many utilities or cloud computing environments. The most notable ones, VMWare, Xen, and KVM, are outlined in the following sections.

VMWare ESXi – VMware is a pioneer in the virtualization market. Its ecosystem of tools ranges from server and desktop virtualization to high-level management tools. ESXi is a VMM (Virtual Machine Manager) from VMware. It is a bare-metal hypervisor, meaning that it installs directly on the physical server, whereas others may require a host operating system. It provides advanced virtualization techniques of processor, memory, and I/O. Especially, through memory ballooning and page sharing, it can overcommit memory, thus increasing the density of VMs inside a single physical server.

Xen—The Xen hypervisor started as an open-source project and has served as a base to other virtualization products, both commercial and open-source. It has pioneered the para-virtualization concept, on which the guest operating system, by means of a specialized kernel, can interact with the hypervisor, thus significantly improving performance. In addition to an

open-source distribution, Xen currently forms the base of commercial hypervisors of several vendors, most notably Citrix XenServer and Oracle VM.

KVM—The kernel-based virtual machine (KVM) is a Linux virtualization subsystem. It has been part of the mainline Linux kernel since version 2.6.20, thus being natively supported by several distributions. In addition, activities such as memory management and scheduling are carried out by existing kernel features, thus making KVM simpler and smaller than hypervisors that take control of the entire machine.

KVM leverages hardware-assisted virtualization, which improves performance and allows it to support unmodified guest operating systems. currently, it supports several versions of Windows, Linux, and UNIX.

4. Virtual Appliances and the Open Virtualization Format

An application combined with the environment needed to run it (operating system, libraries, compilers, databases, application containers, and so forth) is referred to as a —virtual appliance. Packaging application environments in the shape of virtual appliances eases software customization, configuration, and patching and improves portability. Most commonly, an appliance is shaped as a VM disk image associated with hardware requirements, and it can be readily deployed in a hypervisor.

OVF's extensibility has encouraged additions relevant to management of data centers and clouds. Mathews et al. have devised virtual machine contracts (VMC) as an extension to OVF. A VMC aids in communicating and managing the complex expectations that VMs have of their runtime environment and vice versa. A simple example of a VMC is when a cloud consumer wants to specify minimum and maximum amounts of a resource that a VM needs to function. similarly, the cloud provider could express resource limits to bound resource consumption and costs.

5. Autonomic Computing

Autonomic or self-managing, systems rely on monitoring probes and gauges (sensors), on an adaptation engine (autonomic manager) for computing optimizations based on monitoring data, and on effectors to carry out changes on the system. IBM's Autonomic Computing Initiative has contributed to define the four properties of autonomic systems: self-configuration, self- optimization, self-healing, and self-protection. IBM has also suggested a reference model for autonomic control loops of autonomic managers, called MAPE-K (Monitor Analyze Plan Execute—Knowledge)

LAYERS AND TYPES OF CLOUDS

Cloud computing services are divided into three classes, according to the abstraction level of the capability provided and the service model of providers, namely:

1. Infrastructure as a Service,
2. Platform as a Service, and Software as a Service.

Infrastructure as a Service

A cloud infrastructure enables on-demand provisioning of servers running several choices of operating systems and a customized software stack. Infrastructure services are considered as the bottom layer of cloud computing systems. Offering virtualized resources (computation, storage, and communication) on demand is known as Infrastructure as a Service (IaaS).

One of the best examples is Amazon Web Services mainly offers IaaS, which in the case of its EC2 service means offering VMs with a software stack that can be customized similar to how an ordinary physical server would be customized.

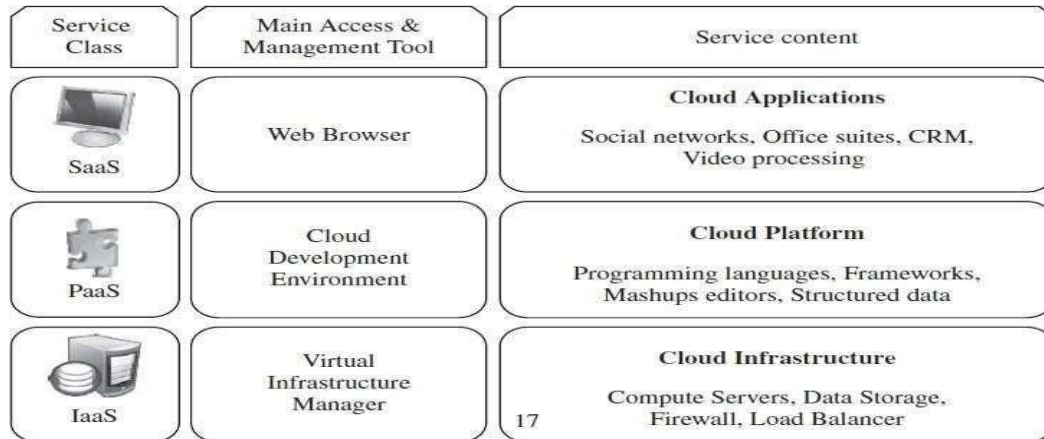


FIGURE 1.3. The cloud computing stack.

Platform as a Service

A *cloud platform* offers an environment on which developers create and deploy applications and do not necessarily need to know how many processors or how much memory that applications will be using. In addition, multiple programming models and specialized services (e.g., data access, authentication, and payments) are offered as building blocks to new applications.

Google AppEngine, an example of Platform as a Service, offers a scalable environment for developing and hosting Web applications, which should be written in specific programming languages such as Python or Java, and use the services' own proprietary structured object data store.

Software as a Service

Traditional desktop applications such as word processing and spreadsheet can now be accessed as a service in the Web. This model of delivering applications, known as Software as a Service (SaaS), alleviates the burden of software maintenance for customers and simplifies development and testing for providers.

Salesforce.com, which relies on the SaaS model, offers business productivity applications (CRM) that reside completely on their servers, allowing customers to customize and access applications on demand.

Deployment Models

Although cloud computing has emerged mainly from the appearance of public

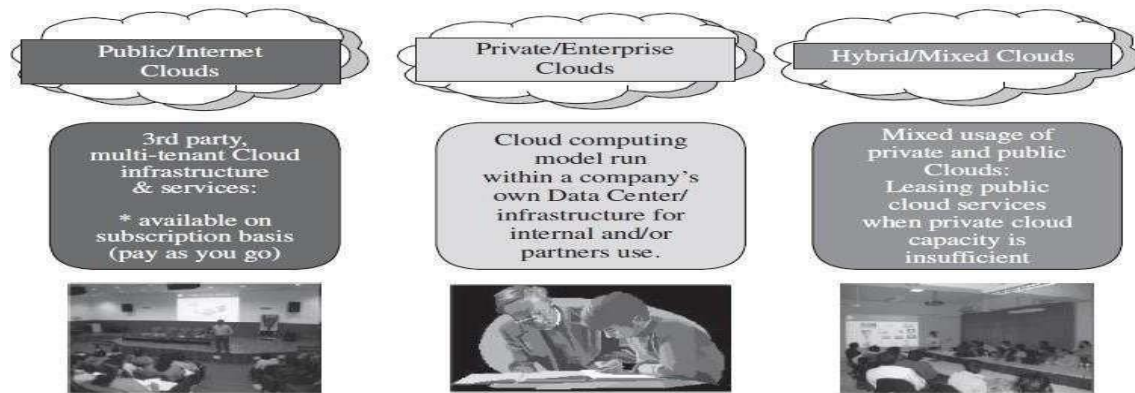


FIGURE 1.4. Types of clouds based on deployment models.

computing utilities, other deployment models, with variations in physical location and distribution, have been adopted. In this sense, regardless of its service class, a cloud can be classified as public, private, community, or hybrid based on model of deployment as shown figure below.

Public cloud & Private cloud: Public cloud as a —cloud made available in a pay-as-you-go manner to the general public. **Private cloud** as —internal data center of a business or other organization, not made available to the general public.

In most cases, establishing a private cloud means restructuring an existing infrastructure by adding virtualization and cloud-like interfaces. This allows users to interact with the local data center while experiencing the same advantages of public clouds, most notably self-service interface, privileged access to virtual servers, and per-usage metering and billing.

A **community cloud** is —shared by several organizations and supports a specific community that has shared concerns (e.g., mission, security requirements, policy, and compliance considerations).

A **hybrid cloud** takes shape when a private cloud is supplemented with computing capacity from public clouds. The approach of temporarily renting capacity to handle spikes in load is known as —cloud-bursting

DESIRED FEATURES OF A CLOUD

Certain features of a cloud are essential to enable services that truly represent the cloud-computing model and satisfy expectations of consumers, and cloud offerings must be having following features:

1. Self-service
2. Per-usage metered and billed
3. Elastic,
4. Customizable.

The following are the features that are explained in detail.

1. Self-Service

Consumers of cloud computing services expect on-demand, nearly instant access to resources. To support this expectation, clouds must allow self-service access so that customers can request, customize, pay, and use services without intervention of human operators.

2. Per-Usage Metering and Billing

Cloud computing eliminates up-front commitment by users, allowing them to request and use only the necessary amount. Services must be priced on a short-term basis (e.g., by the hour), allowing users to release (and not pay for) resources as soon as they are not needed. For these reasons, clouds must implement features to allow efficient trading of service such as pricing, accounting, and billing. Metering should be done accordingly for different types of service (e.g., storage, processing, and bandwidth) and usage promptly reported, thus providing greater transparency.

3. Elasticity

Cloud computing gives the illusion of infinite computing resources available on demand. Therefore, users expect clouds to rapidly provide resources in any quantity at any time. In particular, it is expected that the additional resources can be (a) provisioned, possibly automatically, when an application load increases and (b) released when load decreases (scale up and down).

4. Customization

In a multi-tenant cloud a great disparity between user needs is often the case. Thus, resources rented from the cloud must be highly customizable. In the case of infrastructure services, customization means allowing users to deploy specialized virtual appliances and to be given privileged (root) access to the virtual servers. Other service classes (PaaS and SaaS) offer less flexibility and are not suitable for general-purpose computing, but still are expected to provide a certain level of customization.

CLOUD INFRASTRUCTURE MANAGEMENT

A key challenge IaaS providers face when building a cloud infrastructure is managing physical and virtual resources, namely servers, storage, and networks, in a holistic fashion.

The software toolkit responsible for this orchestration is called a virtual infrastructure manager (VIM). This type of software resembles a traditional operating system—but instead of dealing with a single computer, it aggregates resources from multiple computers, presenting a uniform view to user and applications. Other terms include infrastructure sharing software and virtual infrastructure engine.

There are two categories of tools used to manage cloud they are

5. Cloud toolkits—includes those that —expose a remote and secure interface for creating, controlling and monitoring virtualize resources, but do not specialize in VI management.
6. The virtual infrastructure managers—provide advanced features such as automatic load balancing and server consolidation, but do not expose remote cloud-like interfaces.

The availability of a remote cloud-like interface and the ability of managing many users and their permissions are the primary features that would distinguish —cloud toolkits from —VIMs. However, here we place both categories of tools under the same group (of the VIMs) and, when applicable, we highlight the availability of a remote interface as a feature.

Virtually all VIMs we investigated present a set of basic features related to managing the life cycle of VMs, including networking groups of VMs together and setting up virtual disks for VMs. These basic features pretty much define whether a tool can be used in practical cloud deployments or not. On the other hand, only a handful of software present advanced features (e.g., high availability) which allow them to be used in large-scale production clouds.

Features

We now present a list of both basic and advanced features that are usually available in VIMs.

Virtualization Support— The multi-tenancy aspect of clouds requires multiple customers with disparate requirements to be served by a single hardware infrastructure. Virtualized resources (CPUs, memory, etc.) can be sized and resized with certain flexibility. These features make hardware virtualization, the ideal technology to create a virtual infrastructure that partitions a data center among multiple tenants.

Self-Service, On-Demand Resource Provisioning— Self-service access to resources has been perceived as one the most attractive features of clouds. This feature enables users to directly obtain services from clouds, such as spawning the creation of a server and tailoring its software, configurations, and security policies, without interacting with a human system administrator. This capability eliminates the need for more time-consuming, labor-intensive, human- driven procurement processes familiar to many in IT. Therefore, exposing a self-service interface, through which users can easily interact with the system, is a highly desirable feature of a **Vi manager**.

Multiple Backend Hypervisors— Different virtualization models and tools offer different benefits, drawbacks, and limitations. Thus, some Vi managers provide a uniform management layer regardless of the virtualization technology used.

This characteristic is more visible in open-source Vi managers, which usually provide

pluggable drivers to interact with multiple hypervisors. In this direction, the aim of **libvirt** is to provide a uniform API that Vi managers can use to manage domains (a VM or container running an instance of an operating system) in virtualized nodes using standard operations that abstract hypervisor specific calls.

Storage Virtualization— Virtualizing storage means abstracting logical storage from physical storage. By consolidating all available storage devices in a data center, it allows creating virtual disks independent from device and location. Storage devices are commonly organized in a storage area network (SAN) and attached to servers via protocols such as Fibre Channel, iSCSI, and NFS; a storage controller provides the layer of abstraction between virtual and physical storage.

In the VI management sphere, storage virtualization support is often restricted to commercial products of companies such as VMWare and Citrix.

Interface to Public Clouds— Researchers have perceived that extending the capacity of a local in-house computing infrastructure by borrowing resources from public clouds is advantageous. In this fashion, institutions can make good use of their available resources and, in case of spikes in demand, extra load can be offloaded to rented resources.

A VI manager can be used in a hybrid cloud setup if it offers a driver to manage the life cycle of virtualized resources obtained from external cloud providers. To the applications, the use of leased resources must ideally be transparent.

Virtual Networking— Virtual networks allow creating an isolated network on top of a physical infrastructure independently from physical topology and locations. A virtual LAN (VLAN) allows isolating traffic that shares a switched network, allowing VMs to be grouped into the same broadcast domain. Additionally, a VLAN can be configured to block traffic originated from VMs from other networks. Similarly, the VPN (virtual private network) concept is used to describe a secure and private overlay network on top of a public network (most commonly the public Internet)

Support for creating and configuring virtual networks to group VMs placed throughout a data center is provided by most VI managers. Additionally, VI managers that interface with public clouds often support secure VPNs connecting local and remote VMs.

Dynamic Resource Allocation— In cloud infrastructures, where applications have variable and dynamic needs, capacity management and demand prediction are especially complicated. This fact triggers the need for dynamic resource allocation aiming at obtaining a timely match of supply and demand.

Energy consumption reduction and better management of SLAs can be achieved by dynamically remapping VMs to physical machines at regular intervals. Machines that are not assigned any VM can be turned off or put on a low power state.

A number of VI managers include a dynamic resource allocation feature that continuously monitors utilization across resource pools and reallocates available resources among VMs according to application needs.

Virtual Clusters— Several VI managers can holistically manage groups of VMs. This

feature is useful for provisioning computing virtual clusters on demand, and interconnected VMs for multi-tier Internet applications.

Reservation and Negotiation Mechanism— When users request computational resources to available at a specific time, requests are termed **Advance reservations** (AR), in contrast to best-effort requests, when users request resources whenever available [54]. To support complex requests, such as AR, a VI manager must allow users to —lease resources expressing more complex terms (e.g., the period of time of a reservation). This is especially useful in clouds on which resources are scarce; since not all requests may be satisfied immediately, they can benefit of VM placement strategies that support queues, priorities, and Advance reservations.

Additionally, leases may be negotiated and renegotiated, allowing provider and consumer to modify a lease or present counter proposals until an agreement is reached. This feature is illustrated by the case in which an AR request for a given slot cannot be satisfied, but the provider can offer a distinct slot that is still satisfactory to the user.

High Availability and Data Recovery— The high availability (HA) feature of VI managers aims at minimizing application downtime and preventing business disruption. A few VI managers accomplish this by providing a failover mechanism, which detects failure of both physical and virtual servers and restarts VMs on healthy physical servers. This style of HA protects from host, but not VM, failures.

Data Recovery means data backup in clouds should consider the high data volume involved in VM management. Frequent backup of a large number of VMs, each one with multiple virtual disks attached, should be done with minimal interference in the systems performance. In this sense, some VI managers offer data protection mechanisms that perform incremental backups of VM images. The backup workload is often assigned to proxies, thus offloading production server and reducing network overhead.

Case Studies

In this section, we describe the main features of the most popular VI managers available. Only the most prominent and distinguishing features of each tool are discussed in detail. A detailed side-by-side feature comparison of VI managers is presented in Table 1.1.

Apache VCL: The Virtual Computing project has been incepted in 2004 by researchers at the North Carolina State University as a way to provide customized environments to computer lab users. The software components that support NCSU's initiative have been released as open-source and incorporated by the Apache Foundation.

Since its inception, the main objective of VCL has been providing desktop (virtual lab) and HPC computing environments anytime, in a flexible cost-effective way and with minimal intervention of IT staff. In this sense, VCL was one of the first projects to create a tool with features such as: self-service Web portal, to reduce administrative burden; advance reservation of capacity, to provide resources during classes; and deployment of customized machine images on multiple computers, to provide clusters on demand.

In summary, Apache VCL provides the following features: (i) multi-platform controller, based on Apache/PHP; (ii) Web portal and XML-RPC interfaces; (iii) support for VMware hypervisors (ESX, ESXi, and Server); (iv) virtual networks; (v) virtual clusters; and (vi) advance reservation of capacity.

AppLogic: AppLogic is a commercial VI manager, the flagship product of 3tera Inc. from California, USA. The company has labeled this product as a Grid Operating System.

AppLogic provides a fabric to manage clusters of virtualized servers, focusing on managing multi-tier Web applications. It views an entire application as a collection of components that must be managed as a single entity. Several components such as firewalls, load balancers, Web servers, application servers, and database servers can be set up and linked together. Whenever the application is started, the system manufactures and assembles the virtual infrastructure required to run it. Once the application is stopped, AppLogic tears down the infrastructure built for it.

AppLogic offers dynamic appliances to add functionality such as Disaster Recovery and Power optimization to applications. The key differential of this approach is that additional functionalities are implemented as another pluggable appliance instead of being added as a core functionality of the VI manager.

In summary, 3tera AppLogic provides the following features: Linux-based controller; CLI and GUI interfaces; Xen backend; Global Volume Store (GVS) storage virtualization; virtual networks; virtual clusters; dynamic resource allocation; high availability; and data protection.

Citrix Essentials: The Citrix Essentials suite is one the most feature complete VI management software available, focusing on management and automation of data centers. It is essentially a hypervisor-agnostic solution, currently supporting Citrix XenServer and Microsoft Hyper-V. Citrix Essentials provides the following features: Windows-based controller; GUI, CLI, Web portal, and XML-RPC interfaces which support for XenServer and Hyper-V hypervisors. Citrix Storage Link storage virtualization; virtual networks; dynamic resource allocation; three-level high availability (i.e., recovery by VM restart, recovery by activating paused duplicate VM, and running duplicate VM continuously) [58]; data protection with Citrix Consolidated Backup.

Enomaly ECP: The Enomaly Elastic Computing Platform, in its most complete edition, offers most features a service provider needs to build an IaaS cloud.

Most notably, ECP Service Provider Edition offers a Web-based customer dashboard that allows users to fully control the life cycle of VMs. Usage accounting is performed in real time and can be viewed by users. Similar to the functionality of virtual appliance marketplaces, ECP allows providers and users to package and exchange applications.

In summary, Enomaly ECP provides the following features: Linux-based controller; Web portal and Web services (REST) interfaces; Xen back-end; interface to the Amazon EC2 public cloud; virtual networks; virtual clusters (ElasticValet).

Eucalyptus: The Eucalyptus framework was one of the first open-source projects to focus on building IaaS clouds. It has been developed with the intent of providing an open-source

implementation nearly identical in functionality to Amazon Web Services APIs. Therefore, users can interact with a Eucalyptus cloud using the same tools they use to access Amazon EC2. It also distinguishes itself from other tools because it provides a storage cloud API—emulating the Amazon S3 API—for storing general user data and VM images.

In summary, Eucalyptus provides the following features: Linux-based controller with administration Web portal; EC2-compatible (SOAP, Query) and S3-compatible (SOAP, REST) CLI and Web portal interfaces; Xen, KVM, and VMWare backends; Amazon EBS-compatible virtual storage devices; interface to the Amazon EC2 public cloud; virtual networks.

Nimbus3: The Nimbus toolkit is built on top of the Globus framework. Nimbus provides most features in common with other open-source VI managers, such as an EC2-compatible front-end API, support to Xen, and a backend interface to Amazon EC2. However, it distinguishes from others by providing a Globus Web Services Resource Framework (WSRF) interface. It also provides a backend service, named Pilot, which spawns VMs on clusters managed by a local resource manager (LRM) such as PBS and SGE.

Nimbus' core was engineered around the Spring framework to be easily extensible, thus allowing several internal components to be replaced and also eases the integration with other systems.

In summary, Nimbus provides the following features: Linux-based controller; EC2-compatible (SOAP) and WSRF interfaces; Xen and KVM backend and a Pilot program to spawn VMs through an LRM; interface to the Amazon EC2 public cloud; virtual networks; one-click virtual clusters.

OpenNebula: OpenNebula is one of the most feature-rich open-source VI managers. It was initially conceived to manage local virtual infrastructure, but has also included remote interfaces that make it viable to build public clouds. Altogether, four programming APIs are available: XML-RPC and libvirt for local interaction; a subset of EC2 (Query) APIs and the OpenNebula Cloud API (OCA) for public access.

OpenNebula provides the following features: Linux-based controller; CLI, XML-RPC, EC2-compatible Query and OCA interfaces; Xen, KVM, and VMware backend; interface to public clouds (Amazon EC2, ElasticHosts); virtual networks; dynamic resource allocation; advance reservation of capacity.

OpenPEX: OpenPEX (Open Provisioning and EXecution Environment) was constructed around the notion of using advance reservations as the primary method for allocating VM instances.

OpenPEX provides the following features: multi-platform (Java) controller; Web portal and Web services (REST) interfaces; Citrix XenServer backend; advance reservation of capacity with negotiation.

oVirt: oVirt is an open-source VI manager, sponsored by Red Hat's Emergent Technology group. It provides most of the basic features of other VI managers.

oVirt provides the following features: Fedora Linux-based controller packaged as a virtual appliance; Web portal interface; KVM backend.

Platform ISF: Infrastructure Sharing Facility (ISF) is the VI manager offering from

Platform Computing. The company, mainly through its LSF family of products, has been serving the HPC market for several years.

ISF provides the following features: Linux-based controller packaged as a virtual appliance; Web portal interface; dynamic resource allocation; advance reservation of capacity; high availability.

VMWare vSphere and vCloud: vSphere is VMware's suite of tools aimed at transforming IT infrastructures into private clouds. In the vSphere architecture, servers run on the ESXi platform. A separate server runs vCenter Server, which centralizes control over the entire virtual infrastructure. Through the vSphere Client software, administrators connect to vCenter Server to perform various tasks.

The Distributed Resource Scheduler (DRS) makes allocation decisions based on predefined rules and policies. It continuously monitors the amount of resources available to VMs and, if necessary, makes allocation changes to meet VM requirements. In the storage virtualization realm, vStorage VMFS is a cluster file system to provide aggregate several disks in a single volume. VMFS is especially optimized to store VM images and virtual disks. It supports storage equipment that use Fibre Channel or iSCSI SAN.

vSphere provides the following features: Windows-based controller (vCenter Server); CLI, GUI, Web portal, and Web services interfaces; VMware ESX, ESXi backend; VMware vStorage VMFS storage virtualization; interface to external clouds (VMware vCloud partners); virtual networks (VMware Distributed Switch); dynamic resource allocation (VMware DRM); high availability; data protection (VMware Consolidated Backup).

INFRASTRUCTURE AS A SERVICE PROVIDERS

Public Infrastructure as a Service providers commonly offer virtual servers containing one or more CPUs, running several choices of operating systems and a customized software stack. In addition, storage space and communication facilities are often provided.

Features

IAAS offers a set of specialized features that can influence the cost benefit ratio to be experienced by user applications when moved to the cloud.

The most relevant features are:

1. Geographic distribution of data centers.
2. Variety of user interfaces and APIs to access the system.
3. Specialized components and services that aid Particular applications (e.g., load- balancers, firewalls).
4. Choice of virtualization platform and operating systems and
5. Different billing methods and period (e.g., prepaid vs. postpaid, hourly vs. monthly).

Geographic Presence— To improve availability and responsiveness, a provider of worldwide services would typically build several data centers distributed around the world. For example, Amazon Web Services presents the concept of availability zones and regions

for its EC2 service. Availability zones are distinct locations that are engineered to be insulated from failures in other availability zones and provide inexpensive, low-latency network connectivity to other availability zones in the same region. Regions, in turn, are geographically dispersed and will be in separate geographic areas or countries.

User Interfaces and Access to Servers— Ideally, a public IaaS provider must provide multiple access means to its cloud, thus catering for various users and their preferences. Different types of user interfaces (UI) provide different levels of abstraction, the most common being graphical user interfaces (GUI), command-line tools (CLI), and Web service (WS) APIs.

GUIs are preferred by end users who need to launch, customize, and monitor a few virtual servers and do not necessarily need to repeat the process several times. On the other hand, CLIs offer more flexibility and the possibility of automating repetitive tasks via scripts (e.g., start and shutdown a number of virtual servers at regular intervals).

Advance Reservation of Capacity— Advance reservations allow users to request for an IaaS provider to reserve resources for a specific time frame in the future, thus ensuring that cloud resources will be available at that time. However, most clouds only support best-effort requests that means users can request server whenever resources are available .

Amazon Reserved Instances is a form of advance reservation of capacity, allowing users to pay a fixed amount of money in advance to guarantee resource availability at anytime during an agreed period and then paying a discounted hourly rate when resources are in use. However, only long periods of 1 to 3 years are offered; therefore, users cannot express their reservations in finer granularities—for example, hours or days.

Automatic Scaling and Load Balancing— Automatic scaling is a highly desirable feature of IaaS clouds. It allows users to set conditions for when they want their applications to scale up and down, based on application-specific metrics such as transactions per second, number of simultaneous users, request latency, and so forth.

When the number of virtual servers is increased by automatic scaling, incoming traffic must be automatically distributed among the available servers. This activity enables applications to promptly respond to traffic increase while also achieving greater fault tolerance.

Service-Level Agreement. Service-level agreements (SLAs) are offered by IaaS providers to express their commitment to delivery of a certain QoS. To customers it serves as a warranty. An SLA usually include availability and performance guarantees. Additionally, metrics must be agreed upon by all parties as well as penalties for violating these expectations.

Most IaaS providers focus their SLA terms on availability guarantees, specifying the minimum percentage of time the system will be available during a certain period. For instance, Amazon EC2 states that —if the annual uptime Percentage for a customer drops below 99.95% for the service year, that customer is eligible to receive a service credit equal to 10% of their bill.³¹

Hypervisor and Operating System Choice— Traditionally, IaaS offerings have been based on heavily customized open-source Xen deployments. IaaS providers needed expertise in Linux, networking, virtualization, metering, resource management, and many

other low-level aspects to successfully deploy and maintain their cloud offerings.

More recently, there has been an emergence of turnkey IaaS platforms such as VMWare VCloud and Citrix Cloud Center (C3) which have lowered the barrier of entry for IaaS competitors, leading to a rapid expansion in the IaaS marketplace.

Case Studies

Amazon Web Services: Amazon WS4 (AWS) is one of the major players in the cloud computing market. It pioneered the introduction of IaaS clouds in 2006. It offers a variety of cloud services, most notably: S3 (storage), EC2 (virtual servers), Cloudfront (content delivery), Cloudfront Streaming (video streaming), Simple DB (structured datastore), RDS (Relational Database), SQS (reliable messaging), and Elastic MapReduce (data processing). The ElasticCompute Cloud (EC2) offers Xen-based virtual servers (instances) that can be instantiated from Amazon Machine Images (AMIs). Instances are available in a variety of sizes, operating systems, architectures, and price. CPU capacity of instances is measured in Amazon Compute Units and, although fixed for each instance, vary among instance types from 1 (small instance) to 20 (high CPU instance). Each instance provides a certain amount of non persistent disk space; a persistence disk service (Elastic Block Storage) allows attaching virtual disks to instances with space up to 1TB. Elasticity can be achieved by combining the Cloud Watch, Auto Scaling and Elastic Load Balancing features, which allow the number of instances to scale up and down automatically based on a set of customizable rules, and traffic to be distributed across available instances. Fixed IP address (Elastic IPs) are not available by default, but can be obtained at an additional cost.

Flexiscale: Flexiscale is a UK-based provider offering services similar in nature to Amazon Web Services. Flexiscale cloud provides the following features: available in UK; Web services (SOAP), Web-based user interfaces; access to virtual server mainly via SSH (Linux) and Remote Desktop (Windows); 100% availability SLA with automatic recovery of VMs in case of hardware failure; per hour pricing; Linux and Windows operating systems; automatic scaling (horizontal/vertical).

Joyent: Joyent's Public Cloud offers servers based on Solaris containers virtualization technology. These servers, dubbed accelerators, allow deploying various specialized software- stack based on a customized version of Open- Solaris operating system, which include by default a Web-based configuration tool and several pre-installed software, such as Apache, MySQL, PHP, Ruby on Rails, and Java. Software load balancing is available as an accelerator in addition to hardware load balancers. A notable feature of Joyent's virtual servers is automatic vertical scaling of CPU cores, which means a virtual server can make use of additional CPUs automatically up to the maximum number of cores available in the physical host.

The Joyent public cloud offers the following features: multiple geographic locations in the United States; Web-based user interface; access to virtual server via SSH and Web-based administration tool; 100% availability SLA; per month pricing; OS-level virtualization Solaris containers; Open- Solaris operating systems; automatic scaling(vertical).

GoGrid: GoGrid, like many other IaaS providers, allows its customers to utilize a range of pre- made Windows and Linux images, in a range of fixed instance sizes. GoGrid also offers —value- addedll stacks on top for applications such as high- volume Web serving, e-

Commerce, and database stores. It offers some notable features, such as a —hybrid hosting facility, which combines traditional dedicated hosts with auto-scaling cloud server infrastructure. As part of its core IaaS offerings, GoGrid also provides free hardware load balancing, auto-scaling capabilities, and persistent storage, features that typically add an additional cost for most other IaaS providers.

Rackspace Cloud Servers: Rackspace Cloud Servers is an IaaS solution that provides fixed size instances in the cloud. Cloud Servers offers a range of Linux- based pre-made images. A user can request different-sized images, where the size is measured by requested RAM, not CPU.

PLATFORM AS A SERVICE PROVIDERS

Public Platform as a Service providers commonly offer a development and deployment environment that allow users to create and run their applications with little or no concern to low- level details of the platform. In addition, specific programming languages and frameworks are made available in the platform, as well as other services such as persistent data storage and in memory caches.

Features

Programming Models, Languages, and Frameworks: Programming models made available by IaaS providers define how users can express their applications using higher levels of abstraction and efficiently run them on the cloud platform.

Each model aims at efficiently solving a particular problem. In the cloud computing domain, the most common activities that require specialized models are: processing of large dataset in clusters of computers (MapReduce model), development of request-based Web services and applications; definition and orchestration of business processes in the form of workflows (Workflow model); and high-performance distributed execution of various computational tasks.

For user convenience, PaaS providers usually support multiple programming languages. Most commonly used languages in platforms include Python and Java (e.g., Google AppEngine), .NET languages (e.g., Microsoft Azure), and Ruby (e.g., Heroku). Force.com has devised its own programming language (Apex) and an Excel-like query language, which provide higher levels of abstraction to key platform functionalities.

A variety of software frameworks are usually made available to PaaS developers, depending on application focus. Providers that focus on Web and enterprise application hosting offer popular frameworks such as Ruby on Rails, Spring, Java EE, and .NET.

Persistence Options: A persistence layer is essential to allow applications to record their state and recover it in case of crashes, as well as to store user data. Web and enterprise application developers have chosen relational databases as the preferred persistence method. These databases offer fast and reliable structured data storage and transaction processing, but may lack scalability to handle several peta bytes of data stored in commodity computers. In the cloud computing domain, distributed storage technologies have emerged, which seek to be robust and highly scalable, at the expense of relational structure and convenient query languages.

CASE STUDIES

Aneka: Aneka is a .NET-based service-oriented resource management and development platform. Each server in an Aneka deployment (dubbed Aneka cloud node) hosts the Aneka container, which provides the base infrastructure that consists of services for persistence, security (authorization, authentication and auditing), and communication (message handling and dispatching). Cloud nodes can be either physical server, virtual machines (Xen Server and VMware are supported), and instances rented from Amazon EC2. The Aneka container can also host any number of optional services that can be added by developers to augment the capabilities of an Aneka Cloud node, thus providing a single, extensible framework for orchestrating various application models.

Several programming models are supported by such task models to enable execution of legacy HPC applications and Map Reduce, which enables a variety of data-mining and search applications. Users request resources via a client to a reservation services manager of the Aneka master node, which manages all cloud nodes and contains scheduling service to distribute request to cloud nodes.

App Engine: Google App Engine lets you run your Python and Java Web applications on elastic infrastructure supplied by Google. App Engine allows your applications to scale dynamically as your traffic and data storage requirements increase or decrease. It gives developers a choice between a Python stack and Java. The App Engine serving architecture is notable in that it allows real-time auto- scaling without virtualization for many common types of Web applications. However, such auto-scaling is dependent on the application developer using a limited subset of the native APIs on each platform, and in some instances you need to use specific Google APIs such as URLFetch, Data store, and mem cache in place of certain native API calls. For example, a deployed App Engine application cannot write to the file system directly (you must use the Google Data store) or open a socket or access another host directly (you must use Google URL fetch service). A Java application cannot create a new Thread either.

Microsoft Azure: Microsoft Azure Cloud Services offers developers a hosted .NET Stack (C#, VB.Net, ASP.NET). In addition, a Java & Ruby SDK for .NET Services is also available. The Azure system consists of a number of elements. The Windows Azure Fabric Controller provides auto-scaling and reliability, and it manages memory resources and load balancing. The .NET Service Bus registers and connects applications together. The .NET Access Control identity providers include enterprise directories and Windows LiveID. Finally, the .NET Workflow allows construction and execution of workflow instances.

Force.com: In conjunction with the Salesforce.com service, the Force.com PaaS allows developers to create add-on functionality that integrates into main Salesforce CRM SaaS application. Force.com offers developers two approaches to create applications that can be deployed on its SaaS platform: a hosted Apex or Visualforce application. Apex is a proprietary Java-like language that can be used to create Salesforce applications. Visual force is an XML-like syntax for building UIs in HTML, AJAX, or Flex to overlay over the Salesforce hosted CRM system. An application store called App Exchange is also provided, which offers a paid & free application directory.

Heroku: Heroku is a platform for instant deployment of Ruby on Rails Web applications. In the Heroku system, servers are invisibly managed by the platform and are never exposed to users. Applications are automatically dispersed across different CPU cores and servers,

maximizing performance and minimizing contention. Heroku has an advanced logic layer that can automatically route around failures, ensuring seamless and uninterrupted service at all times.

CHALLENGES AND RISKS

Despite the initial success and popularity of the cloud computing paradigm and the extensive availability of providers and tools, a significant number of challenges and risks are inherent to this new model of computing. Providers, developers, and end users must consider these challenges and risks to take good advantage of cloud computing. Issues to be faced include user privacy, data security, data lock-in, availability of service, disaster recovery, performance, scalability, energy- efficiency, and programmability.

Security, Privacy, and Trust: Security and privacy affect the entire cloud computing stack, since there is a massive use of third-party services and infrastructures that are used to host important data or to perform critical operations. In this scenario, the trust toward providers is fundamental to ensure the desired level of privacy for applications hosted in the cloud. Legal and regulatory issues also need attention. When data are moved into the Cloud, providers may choose to locate them anywhere on the planet. The physical location of data centers determines the set of laws that can be applied to the management of data. For example, specific cryptography techniques could not be used because they are not allowed in some countries. Similarly, country laws can impose that sensitive data, such as patient health records, are to be stored within national borders.

Data Lock-In and Standardization: A major concern of cloud computing users is about having their data locked-in by a certain provider. Users may want to move data and applications out from a provider that does not meet their requirements. However, in their current form, cloud computing infrastructures and platforms do not employ standard methods of storing user data and applications. Consequently, they do not interoperate and user data are not portable.

The answer to this concern is standardization. In this direction, there are efforts to create open standards for cloud computing. The Cloud Computing Interoperability Forum (CCIF) was formed by organizations such as Intel, Sun, and Cisco in order to —enable a global cloud computing ecosystem whereby organizations are able to seamlessly work together for the purposes for wider industry adoption of cloud computing technology.¶ The development of the Unified Cloud Interface (UCI) by CCIF aims at creating a standard programmatic point of access to an entire cloud infrastructure. In the hardware virtualization sphere, the Open Virtual Format (OVF) aims at facilitating packing and distribution of software to be run on VMs so that virtual appliances can be made portable—that is, seamlessly run on hypervisor of different vendors.

Availability, Fault-Tolerance, and Disaster Recovery: It is expected that users will have certain expectations about the service level to be provided once their applications are moved to the cloud. These expectations include availability of the service, its overall performance, and what measures are to be taken when something goes wrong in the system or its components. In summary, users seek for a warranty before they can comfortably move their business to the cloud. SLAs, which include QoS requirements, must be ideally set up between customers and cloud computing providers to act as warranty. An SLA specifies the details of the service to be provided, including availability and performance

guarantees.

Additionally, metrics must be agreed upon by all parties, and penalties for violating the expectations must also be approved.

Resource Management and Energy-Efficiency: One important challenge faced by providers of cloud computing services is the efficient management of virtualized resource pools. Physical resources such as CPU cores, disk space, and network bandwidth must be sliced and shared among virtual machines running potentially heterogeneous workloads. The multi-dimensional nature of virtual machines complicates the activity of finding a good mapping of VMs onto available physical hosts while maximizing user utility. Dimensions to be considered include: number of CPUs, amount of memory, size of virtual disks, and network bandwidth. Dynamic VM mapping policies may leverage the ability to suspend, migrate, and resume VMs as an easy way of preempting low-priority allocations in favor of higher-priority ones. Migration of VMs also brings additional challenges such as detecting when to initiate a migration, which VM to migrate, and where to migrate. In addition, policies may take advantage of live migration of virtual machines to relocate data center load without significantly disrupting running services. In this case, an additional concern is the trade-off between the negative impact of a live migration on the performance and stability of a service and the benefits to be achieved with that migration. Another challenge concerns the outstanding amount of data to be managed in various VM management activities. Such data amount is a result of particular abilities of virtual machines, including the ability of traveling through space (i.e., migration) and time (i.e., check pointing and rewinding), operations that may be required in load balancing, backup, and recovery scenarios. In addition, dynamic provisioning of new VMs and replicating existing VMs require efficient mechanisms to make VM block storage devices (e.g., image files) quickly available at selected hosts. Data centers consumer large amounts of electricity. According to a data published by HP, 100 server racks can consume 1.3MW of power and another 1.3 MW are required by the cooling system, thus costing USD 2.6 million per year. Besides the monetary cost, data centers significantly impact the environment in terms of CO2 emissions from the cooling systems.

COMPUTING PARADIGMS:

High-Performance Computing:

For many years, HPC systems emphasize the raw speed performance. The speed of HPC systems has increased from Gflops in the early 1990s to now Pflops in 2010. This improvement was driven mainly by the demands from scientific, engineering, and manufacturing communities.

Top 500 most powerful computer systems in the world are measured by floating-point speed in Linpack benchmark results. However, the number of supercomputer users is limited to less than 10% of all computer users.

Today, the majority of computer users are using desktop computers or large servers when they conduct Internet searches and market-driven computing tasks.

Three New Computing Paradigms

- With the introduction of SOA, **Web 2.0 services** become available.

- Advances in virtualization make it possible to see the growth of **Internet clouds** as a new computing paradigm.
- The maturity of radio-frequency identification (RFID), Global Positioning System (GPS), and sensor technologies has triggered the development of the **Internet of Things (IoT)**.

Computing Paradigm Distinctions

- Centralized computing
- Parallel computing
- Distributed computing
- Cluster computing
- Cloud computing
- Grid computing

Centralized computing:

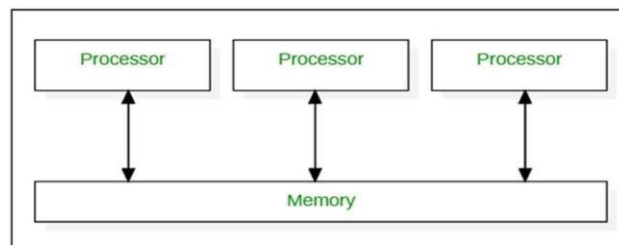
This is a computing paradigm by which all computer **resources are centralized in one physical system.**

All resources (processors, memory, and storage) are fully shared and tightly coupled within one integrated OS. Many data centers and supercomputers are centralized systems, but they are used in parallel, distributed, and cloud computing applications.

Parallel Computing:

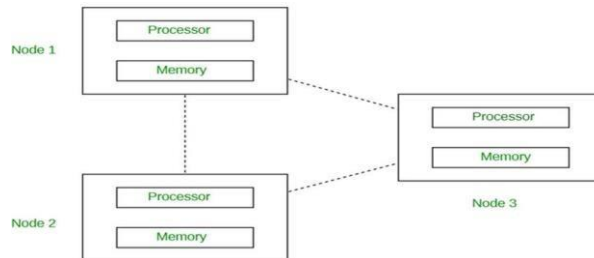
Parallel computing is defined as a type of computing where multiple computer systems are used simultaneously. Here a problem is broken into sub-problems and then further broken down into instructions. These instructions from each sub-problem are executed concurrently on different processors.

Here in the below diagram, you can see how the parallel computing system consists of multiple processors that communicate with each other and perform multiple tasks over a shared memory simultaneously.



Distributed Computing:

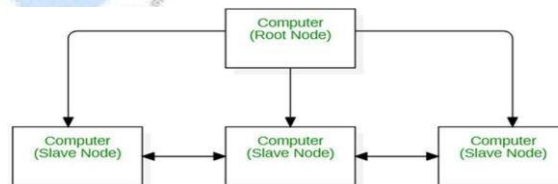
Distributed computing is defined as a type of computing where multiple computer systems work on a single problem. Here all the computer systems are linked together, and the problem is divided into sub-problems where each part is solved by different computer systems. The goal of distributed computing is to increase the performance and efficiency of the system and ensure fault tolerance. In the below diagram, each processor has its own local memory, and all the processors communicate with each other over a network.



Cluster Computing:

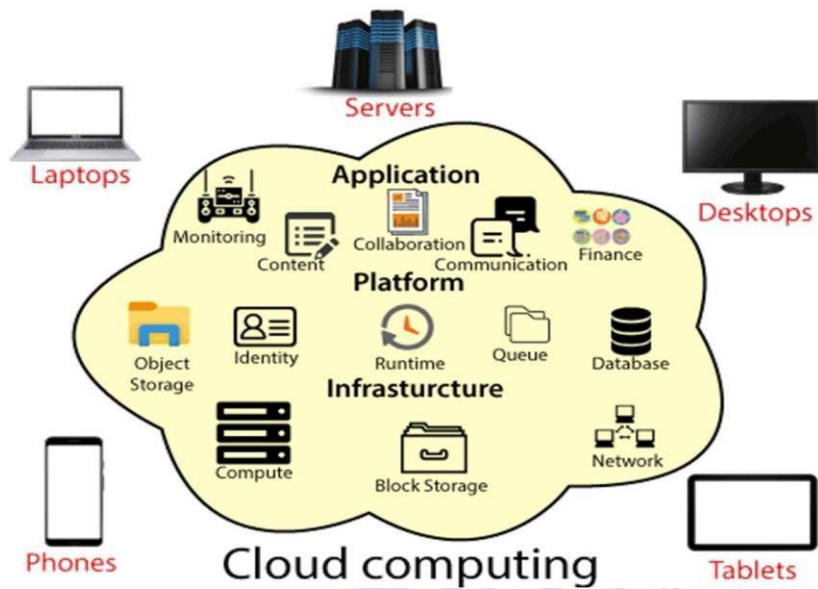
A cluster is a group of independent computers that work together to perform the tasks given. Cluster computing is defined as a type of computing that consists of two or more independent computers, referred to as nodes, that work together to execute tasks as a single machine.

The goal of cluster computing is to increase the performance, scalability and simplicity of the system. As you can see in the below diagram, all the nodes, (irrespective of whether they are a parent node or child node), act as a single entity to perform the tasks.



Cloud computing :

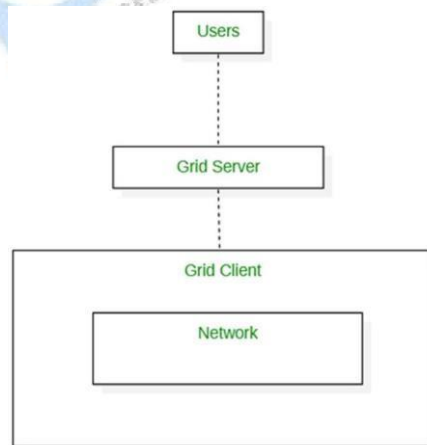
An Internet cloud of resources can be either a centralized or a distributed computing system. The cloud applies parallel or distributed computing, or both. Clouds can be built with physical or virtualized resources over large data centers that are centralized or distributed.



Grid Computing:

Grid computing is defined as a type of computing where it constitutes a network of computers that work together to perform tasks that may be difficult for a single machine to handle. All the computers on that network work under the same umbrella and are termed as a virtual supercomputer.

The tasks they work on are of either high computing power and consist of large data sets. All communication between the computer systems in grid computing is done on the —data grid. The goal of grid computing is to solve more high computational problems in less time and improve productivity.



UNIT - 2

MIGRATING INTO A CLOUD: INTRODUCTION

Cloud computing has been a hotly debated and discussed topic amongst IT professionals and researchers both in the industry and in academia. There are intense discussions on several blogs, in Web sites, and in several research efforts. This also resulted in several entrepreneurial efforts to help leverage and migrate into the cloud given the myriad issues, challenges, benefits, and limitations and lack of comprehensive understanding of what cloud computing can do.

On the one hand, there were these large cloud computing IT vendors like Google, Amazon, and Microsoft, who had started offering cloud computing services on what seemed like a demonstration and trial basis though not explicitly mentioned. They were charging users fees that in certain contexts demonstrated very attractive pricing models.

Most enterprises today are powered by captive data centers. In most large or small enterprises today, IT is the backbone of their operations. Invariably for these large enterprises, their data centers are distributed across various geographies.

They comprise systems and software that span several generations of products sold by a variety of IT vendors. In order to meet varying loads, most of these data centers are provisioned with capacity beyond the peak loads experienced.

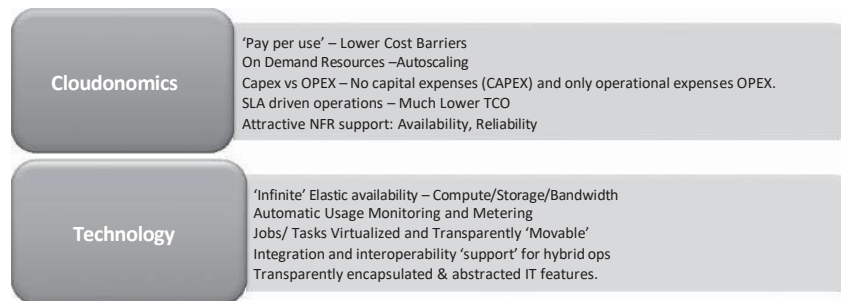
Many data center management teams have been continuously innovating their management practices and technologies.

Cloud computing turned attractive to them because they could pass on the additional demand from their IT setups onto the cloud while paying only for the usage and being unencumbered by the load of operations and management.

The Promise of the Cloud

1. The promise of the cloud both on the business front (the attractive cloudonomics) and the technology front widely aided the CxOs to spawn out several non-mission critical IT needs from the ambit of their captive traditional data centers to the appropriate cloud service.
2. Invariably, these IT needs had some common features: They were typically Web-oriented; they represented seasonal IT demands; they were amenable to parallel batch processing; they were non-mission critical and therefore did not have high security demands. They included scientific applications too [7]. Several small and medium business enterprises, however, leveraged the cloud much beyond the cautious user.
3. Many startups opened their IT departments exclusively using cloud services—very successfully and with high ROI. Having observed these successes, several large enterprises have started successfully running pilots for leveraging the cloud.
4. Many large enterprises run SAP to manage their operations. SAP itself is experimenting with running its suite of products: SAP Business One as well as

- SAP NetWeaver on Amazon cloud offerings.
5. Gartner, Forrester, and other industry research analysts predict that a substantially significant percentage of the top enterprises in the world would have migrated most of their IT needs to the cloud offerings by 2012, thereby demonstrating the widespread impact and benefits from cloud computing. Indeed, the promise of the cloud has been significant in its impact.



The promise of the cloud computing services.

BROAD APPROACHES TO MIGRATING INTO THE CLOUD

Why Migrate:

There are economic and business reasons why an enterprise application can be migrated into the cloud, and there are also several number of technological reasons.

Many of these efforts come up as initiatives in adoption of cloud technologies in the enterprise, resulting in integration of enterprise applications running off the captive data centers with the new ones that have been developed on the cloud.

At the core, migration of an application into the cloud can happen in one of several ways:

1. Either the application is clean and independent.
2. Perhaps some degree of code needs to be modified and adapted or the design (and therefore the code) needs to be first migrated into the cloud computing service environment
3. Perhaps the migration results in the core architecture being migrated for a cloud computing service setting, this resulting in a new architecture being developed, along with the accompanying design and code implementation.
4. Perhaps while the application is migrated as is, it is the usage of the application that needs to be migrated and therefore adapted and modified.
5. Migration can happen at five levels i.e.,
 1. Application
 2. Code
 3. Design
 4. Architecture
 5. Usage

With due simplification, the migration of an enterprise application is best captured by the following:

$$P \rightarrow P_C^0 \cup P_L^0 \rightarrow P_{OFC}^0 \cup P_L^0$$

1. Where P is the application before migration running in captive data center.
2. P_C^0 is the application part after migration into a (hybrid) cloud.
3. P_L^0 is the part of application being run in the captive local data center.
4. P_{OFC}^0 is the application part optimized for cloud.

Seven-Step Model of Migration into a Cloud



Step-1: Cloud migration assessments comprise assessments to understand the issues involved in the specific case of migration at the application level or the code, the design, the architecture, or usage levels. These assessments are about the cost of migration as well as about the ROI that can be achieved in the case of production version.

Step-2: Isolating all systemic and environmental dependencies of the enterprise application components within the captive data center.

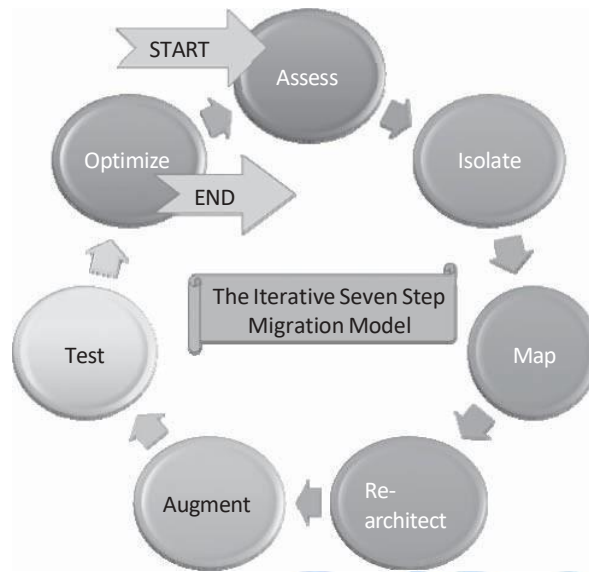
Step-3: Generating the mapping constructs between what shall possibly remain in the local captive data center and what goes onto the cloud.

Step-4: substantial part of the enterprise application needs to be rearchitected, redesigned, and reimplemented on the cloud

Step-5: We leverage the intrinsic features of the cloud computing service to augment our enterprise application in its own small ways.

Step-6: We validate and test the new form of the enterprise application with an extensive test suite that comprises testing the components of the enterprise application on the cloud as well

Step-7: Test results could be positive or mixed. In the latter case, we iterate and optimize as appropriate. After several such optimizing iterations, the migration is deemed successful



Enriching the ‘Integration as a Service’ Paradigm for the Cloud Era

1. The trend-setting cloud paradigm represents the cool conglomeration of a number of proven and promising Web and enterprise technologies. Though the cloud idea is not conceptually new, practically it has brought in myriad tectonic shifts for the whole information and communication technology (ICT) industry.
2. The cloud concepts have progressively and perceptibly impacted the IT and business domains on several critical aspects. The cloud computing has brought in series of novelty-packed deployment, delivery, consumption and pricing models whereas the service orientation prescribes as much simpler application design mechanism.
3. The noteworthy contribution of the much-discoursed and deliberated cloud computing is the faster realization and proliferation of dynamic, converged, adaptive, on-demand, and online compute infrastructures, which are the key requirement for the future IT.
4. The delightful distinctions here are that clouds guarantee most of the nonfunction requirements (Quality of Service (QoS) attributes) such as availability, high performance, on-demand scalability/elasticity, affordability, global-scale accessibility and usability, energy efficiency etc.
5. Having understood the exceptional properties of cloud infrastructures (hereafter will be described as just clouds), most of the global enterprises (small, medium and even large) are steadily moving their IT offerings such as business services and applications to clouds. This transition will facilitate a higher and deeper reach and richness in application delivery and consumability.
6. Product vendors having found that the cloud style is a unique proposition are moving their platforms, databases, and middleware to clouds. Cloud Infrastructure providers

- are establishing cloud centers to host a variety of ICT services and platforms of worldwide individuals, innovators, and institutions.
7. Cloud service providers (CSPs) are very aggressive in experimenting and embracing the cool cloud ideas and today every business and technical services are being hosted in clouds to be delivered to global customers, clients and consumers over the Internet communication infrastructure.

INTEGRATION AS A SERVICE (IAAS):

Why Integration?

- Increasingly business applications are deployed in clouds to reap the business and technical benefits.
- On the other hand, there are still innumerable applications and data sources locally stationed and sustained primarily due to the security reason.
- The question here is how to create a seamless connectivity between those hosted and on-premise applications to empower them to work together.

How Integration is done?

- Integration as a service (IaaS) is the budding and distinctive capability of clouds in fulfilling the business integration requirements.
- IaaS overcomes these challenges by smartly utilizing the time-tested business- to-business (B2B) integration technology as the value added bridge between SaaS solutions and in-house business applications.

SaaS INTEGRATION

- ◆ Cloud-centric integration solutions are being developed and demonstrated for showcasing their capabilities for integrating enterprise and cloud applications.
- ◆ Now with the arrival and adoption of the transformative and disruptive paradigm of cloud computing, every ICT products are being converted into a collection of services to be delivered via the open Internet
- ◆ In that line, the standards-compliant integration suites are being transitioned into services so that any integration need of any one from any part of the world, can be easily, cheaply and rapidly met.

Integration as a Service (IaaS):

Migration of the functionality of a typical enterprise application integration (EAI) hub / enterprise service bus (ESB) into the cloud for providing for smooth data transport between any enterprise and SaaS applications.

- Users subscribe to IaaS as they would do for any other SaaS application.
- Cloud middleware will be made available as a service.
- For service integration, it is enterprise service bus (ESB) and for data integration, it is enterprise data bus (EDB).
- There are Message oriented middleware (MOM) and message brokers for integrating decoupled applications through message passing and pickup.
- Events are coming up fast and there are complex event processing (CEP) engines that receive a stream of diverse events from diverse sources, process them at real- time to

extract and figure out the encapsulated knowledge, and accordingly select and activate one or more target applications.

- Cloud infrastructure is not very useful without SaaS applications that run on top of them, and SaaS applications are not very valuable without access to the critical corporate data that is typically locked away in various corporate systems.
- Cloud applications to offer maximum value to their users, they need to provide a simple mechanism to import or load external data, export or replicate their data for reporting or analysis purposes, and finally keep their data synchronized with on-premise applications.

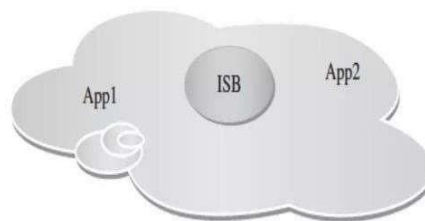
Why SaaS Integration is hard? Reasons:

- **Limited Access:** Access to cloud resources (SaaS, PaaS, and the infrastructures) is more limited than local applications. Once applications move to the cloud, custom applications must be designed to support integration because there is no longer that low level of access. Enterprises putting their applications in the cloud or those subscribers of cloud-based business services are dependent on the vendor to provide the integration hooks and APIs.
- **Dynamic Resources:** Cloud resources are virtualized and service-oriented. That is, everything is expressed and exposed as a service. Due to the dynamism factor infrastructural changes are liable for dynamic changes. These would clearly impact the integration model.
- **Performance:** Clouds support application scalability and resource elasticity. However the network distances between elements in the cloud are no longer under our control. Because of the round trip latency, the cloud integration performance is bound to slow down

NEW INTEGRATION SCENARIOS

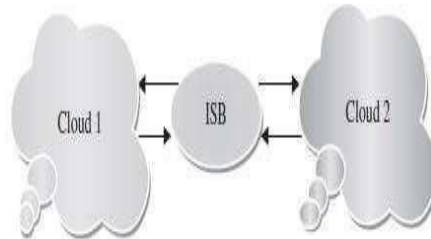
Three major integration scenarios

Within a Public Cloud: Two different applications are hosted in a cloud. The role of the cloud integration middleware (say cloud-based ESB or internet service bus (ISB)) is to seamlessly enable these applications to talk to each other. These applications can be owned by two different companies. They may live in a single physical server but run on different virtual machines.

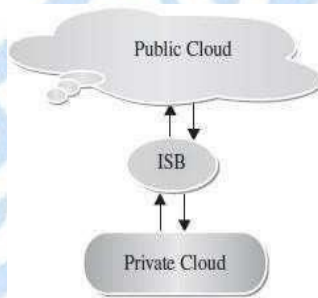


Homogeneous Clouds: The applications to be integrated are positioned in two geographically

separated cloud infrastructures. The integration middleware can be in cloud 1 or 2 or in a separate cloud. There is a need for data and protocol transformation and they get done by the ISB



Heterogeneous Clouds : One application is in public cloud and the other application is in private cloud



The Integration Methodologies

There are three types for cloud integration

- **Traditional Enterprise Integration Tools can be empowered with special connectors to access Cloud-located Applications:** With a persistent rise in the necessity towards accessing and integrating cloud applications, special drivers, connectors and adapters are being built and incorporated on the existing integration platforms to enable bidirectional connectivity with the participating cloud services.
- **Traditional Enterprise Integration Tools are hosted in the Cloud:** This approach is similar to the first option except that the integration softwaresuite is now hosted in any third-party cloud infrastructures so that the enterprise does not worry about procuring and managing the hardware or installing the integration software
- **Integration-as-a-Service (IaaS) or On-Demand Integration Offerings:** These are SaaS applications that are designed to deliver the integration service securely over the Internet and are able to integrate cloud applications with the on-premise systems, cloud-to-cloud applications.

Characteristics of Integration Solutions and Products

- Connectivity refers to the ability of the integration engine to engage with both the source and target systems using available native interfaces
- Semantic Mediation refers to the ability to account for the differences between application semantics between two or more systems. Semantics means how information gets understood, interpreted and represented within information systems
- Data Mediation converts data from a source data format into destination data format. Coupled with semantic mediation, data mediation or data transformation is the process of converting data from one native format on the source system, to another data format for the target system
- Data Migration is the process of transferring data between storage types, formats, or systems. Data migration means that the data in the old system is mapped to the new systems
- Data Security means the ability to insure that information extracted from the source systems has to securely be placed into target systems
- Data Integrity means data is complete and consistent. Thus, integrity has to be guaranteed when data is getting mapped and maintained during integration operations, such as data synchronization between on-premise and SaaS-
- Based systems.
- Governance refers to the processes and technologies that surround a system or systems, which control how those systems are accessed and leveraged

Products And Platforms

- Jitter bit is a fully graphical integration solution that provides users a versatile platform and a suite of productivity tools to reduce the Integration efforts sharply.
- It can be used standalone or with existing EAI infrastructures, enabling users to create new projects or consume and modify existing ones offered by the open source community or service provider.
- The Jitter bit solution enables the cool integration among confidential and corporate data, enterprise applications, web services, XML data sources, legacy systems, simple and complex flat files.

Jitter bit is comprised of two major components:

- **Jitter bit Integration Environment:** An intuitive point-and-click graphical UI that enables to quickly configure, test, deploy and manage integration projects on the Jitterbit server.
- **Jitterbit Integration Server:** A powerful and scalable run-time engine that processes all the integration operations, fully configurable and manageable from the Jitterbit application.
- **Boomi Software:** BoomiAtomSphere is an integration service that is completely on-demand and connects any combination of SaaS, PaaS, cloud, and on-premise applications without the burden of installing and maintaining software packages or appliances. Anyone can securely build, deploy and manage simple to complex integration processes using only web browser.
- **Bungee Connect:** Bungee Connect enables cloud computing by offering an application development and deployment platform that enables highly interactive

applications integrating multiple data sources and facilitating instant deployment.

- Built specifically for cloud development, Bungee Connect reduces the efforts to integrate (mashup) multiple web services into a single application.
- **OpSource Connect:** Expands on the OpSource Services Bus (OSB) by providing the infrastructure for two-way web services interactions, allowing customers to consume and publish applications across a common web services infrastructure.
- **SnapLogic:** SnapLogic is a capable, clean, and uncluttered solution for data integration that can be deployed in enterprise as well as in cloud landscapes. The free community edition can be used for the most common point-to-point data integration tasks

The Pervasive DataCloud: Pervasive Data Cloud is the first multi-tenant platform for delivering the following

- Integration as a Service (IaaS) for both hosted and on-premises applications and data sources
- Packaged turnkey integration
- Integration that supports every integration scenario
- Connectivity to hundreds of different applications and data sources

Other Products

- Bluewolf
- Online MQ
- CloudMQ
- Linxter

ADVANTAGES OF CLOUD INTEGRATION: Cloud integration offers the following advantages over older, compartmentalized organizational methods:

1. Each user can access personal data in real time from any device and from any location with Internet access.
2. Each user can integrate personal data such as calendars and contact lists served by diverse application programs.
3. Each user can employ the same login information (username and password) for all personal applications.
4. The system efficiently passes control messages among application programs.
5. By avoiding the use of data silos, data integrity is maintained.
6. Cloud integration offers scalability to allow for future expansion in terms of the number of users, the number of applications, or both.

THE ONSET OF KNOWLEDGE ERA

1. Path-breaking and people-centric technologies (miniaturization, virtualization, federation, composition, collaboration, etc.) are emerging and are being

experimented, expounded, and established in order to empower the professional and the personal IT to be smart, simple, supple and sensitive towards users' situational needs and to significantly enhance peoples' comfort, care, convenience and choice.

2. In the monolithic mainframe era, one centralized and large system performed millions of operations to respond to thousands of users (one-to-many).
3. Today everyone has his own compute machine (one-to-one), and tomorrow a multitude of smart objects and electronic devices (nomadic, wearable, portable, implantable etc.) will seamlessly and spontaneously co- exist, corroborate, correlate, and coordinate with one another dynamically with dexterity to understand one or more users' needs, conceive, construct, and deliver them at right time at right place (many-to-one). Anytime anywhere computing tends towards everywhere, every time and everything computing.
4. Ambient intelligence (AmI) is the newest buzzword today with ambient sensing, networking, perception, decision-making and actuation technologies. Multimedia and multimodal technologies are flourishing in order to be make human interaction more friendly and fruitful.
5. Dynamic virtualized and autonomic infrastructures, flexible, integrated and lean processes, constructive and contributive building-blocks (service, model, composite, agent, aspect etc.), slim and sleek devices and appliances, smart objects empowered by invisible tags and stickers, natural interfaces, ad-hoc and situational networking capabilities all combine adaptively together to accomplish the grandiose goals of the forthcoming ambient intelligence days and decades. In short, IT- sponsored and splurged smartness in every facet of our living in this world is the vision. Software engineering is on the right track with the maturity of service orientation concepts and software as a service (SaaS) model.
6. Clouds chip in mightily in realizing the much-acclaimed knowledge era. Technologies form a dynamic cluster in real-time in order to contribute immensely and immeasurably for all the existing, evolving and exotic expectations of people.

THE EVOLUTION OF SaaS

1. SaaS paradigm is on fast track due to its innate powers and potentials. Executives, entrepreneurs, and end-users are ecstatic about the tactic as well as strategic success of the emerging and evolving SaaS paradigm.
2. A number of positive and progressive developments started to grip this model. Newer resources and activities are being consistently readied to be delivered as a service.
3. Experts and evangelists are in unison that cloud is to rock the total IT community as the best possible infrastructural solution for effective service delivery. There are several ways clouds can be leveraged inspiring and incredibly for diverse IT

problems. Today there is a small list of services being delivered via the clouds and in future, many more critical applications will be deployed and consumed.

4. In short, clouds are set to decimate all kinds of IT inflexibility and dawn a growing array of innovations to prepare the present day IT for sustainable prosperity.
5. IT as a Service (ITaaS) is the most recent and efficient delivery method in the decisive IT landscape. With the meteoric and mesmerizing rise of the service orientation principles, every single IT resource, activity and infrastructure is being viewed and visualized as a service that sets the tone for the grand unfolding of the dreamt service era.
6. These days, systems are designed and engineered as elegant collections of enterprising and evolving services. Infra- structures are service enabled to be actively participative and collaborative.
7. In the same tenor, the much-maligned delivery aspect too has gone through several transformations and today the whole world has solidly settled for the green paradigm IT as a service (ITaaS).
8. This is accentuated due to the pervasive Internet. Also, we are bombarded with innumerable implementation technologies and methodologies. Clouds, as indicated above, is the most visible and viable infrastructure for realizing ITaaS.
9. Another influential and impressive factor is the maturity obtained in the consumption-based metering and billing capability. HP even proclaims this evolving trend as everything as a service.
10. Integration as a service (IaaS) is the budding and distinctive capability of clouds in fulfilling the business integration requirements. Increasingly business applications are deployed in clouds to reap the business and technical benefits. On the other hand, there are still innumerable applications and data sources locally stationed and sustained primarily due to the security reason. The question here is how to create a seamless connectivity between those hosted and on-premises applications to empower them to work together.
11. IaaS over- comes these challenges by smartly utilizing the time-tested business-to-

business (B2B) integration technology as the value-added bridge between SaaS solutions and in-house business applications.

12. B2B systems are capable of driving this new on-demand integration model because they are traditionally employed to automate business processes between manufacturers and their trading partners. That means they provide application-to-application connectivity along with the functionality that is very crucial for linking internal and external software securely. Unlike the conventional EAI solutions designed only for internal data sharing, B2B platforms have the ability to encrypt files for safe passage across the public network, manage large data volumes, transfer batch files, convert disparate file formats, and guarantee data delivery across multiple enterprises.
 13. IaaS just imitates this established communication and collaboration model to create reliable and durable linkage for ensuring smooth data passage between traditional and cloud systems over the Web infrastructure.
 14. The use of hub & spoke (H&S) architecture further simplifies the implementation and avoids placing an excessive processing burden on the customer sides. The hub is installed at the SaaS provider's cloud center to do the heavy lifting such as reformatting files. A spoke unit at each user site typically acts as basic data transfer utility.
 15. With these pieces in place, SaaS providers can offer integration services under the same subscription / usage-based pricing model as their core offerings. This trend of moving all kinds of common and centralised services to clouds is gaining momentum these days.
 16. As resources are getting distributed and decentralised, linking and leveraging them for multiple purposes need a multifaceted infrastructure.
 17. Clouds, being the Web-based infrastructures are the best fit for hosting scores of unified and utility-like platforms to take care of all sorts of brokering needs among connected and distributed ICT systems.
1. The Web is the largest digital information superhighway

2. The Web is the largest repository of all kinds of resources such as web pages, applications comprising enterprise components, business services, beans, POJOs, blogs, corporate data, etc.
 3. The Web is turning out to be the open, cost-effective and generic business execution platform (E-commerce, business, auction, etc. happen in the web for global users) comprising a wider variety of containers, adaptors, drivers, connectors, etc.
 4. The Web is the global-scale communication infrastructure (VoIP, Video conferencing, IP TV etc.)
 5. The Web is the next-generation discovery, Connectivity, and integration middleware
- Thus, the unprecedented absorption and adoption of the Internet is the key driver for the continued success of the cloud computing.

UNIT – III

Infrastructure as a Service (IAAS) & Platform (PAAS): Virtual machines provisioning and Migration services, Virtual Machines Provisioning and Manageability, Virtual Machine Migration Services, VM Provisioning and Migration in Action. On the Management of Virtual machines for Cloud Infrastructures- Aneka—Integration of Private and Public Clouds.

INFRASTRUCTURE AS A SERVICE (IAAS) & PLATFORM (PAAS)

INFRASTRUCTURE AS A SERVICE PROVIDERS

Public Infrastructure as a Service providers commonly offer virtual servers containing one or more CPUs, running several choices of operating systems and a customized software stack. In addition, storage space and communication facilities are often provided.

Features

IAAS offers a set of specialized features that can influence the cost benefit ratio to be experienced by user applications when moved to the cloud.

The most relevant features are:

1. Geographic distribution of data centers.
2. Variety of user interfaces and APIs to access the system.
3. Specialized components and services that aid Particular applications (e.g., load- balancers, firewalls).
4. Choice of virtualization platform and operating systems and
5. Different billing methods and period (e.g., prepaid vs. postpaid, hourly vs. monthly).

Geographic Presence— To improve availability and responsiveness, a provider of worldwide services would typically build several data centers distributed around the world. For example, Amazon Web Services presents the concept of availability zones and regions for its EC2 service. Availability zones are distinct locations that are engineered to be insulated from failures in other availability zones and provide inexpensive, low-latency network connectivity to other availability zones in the same region. Regions, in turn, are geographically dispersed and will be in separate geographic areas or countries.

User Interfaces and Access to Servers— Ideally, a public IaaS provider must provide multiple access means to its cloud, thus catering for various users and their preferences. Different types of user interfaces (UI) provide different levels of abstraction, the most common being graphical user interfaces (GUI), command-line tools (CLI), and Web service (WS) APIs.

GUIs are preferred by end users who need to launch, customize, and monitor a few virtual servers and do not necessarily need to repeat the process several times. On the other hand, CLIs offer more flexibility and the possibility of automating repetitive tasks via

scripts (e.g., start and shutdown a number of virtual servers at regular intervals).

Advance Reservation of Capacity— Advance reservations allow users to request for an IaaS provider to reserve resources for a specific time frame in the future, thus ensuring that cloud resources will be available at that time. However, most clouds only support best-effort requests that means users can request server whenever resources are available .

Amazon Reserved Instances is a form of advance reservation of capacity, allowing users to pay a fixed amount of money in advance to guarantee resource availability at anytime during an agreed period and then paying a discounted hourly rate when resources are in use. However, only long periods of 1 to 3 years are offered; therefore, users cannot express their reservations in finer granularities—for example, hours or days.

Automatic Scaling and Load Balancing— Automatic scaling is a highly desirable feature of IaaS clouds. It allows users to set conditions for when they want their applications to scale up and down, based on application-specific metrics such as transactions per second, number of simultaneous users, request latency, and so forth.

When the number of virtual servers is increased by automatic scaling, incoming traffic must be automatically distributed among the available servers. This activity enables applications to promptly respond to traffic increase while also achieving greater fault tolerance.

Service-Level Agreement- Service-level agreements (SLAs) are offered by IaaS providers to express their commitment to delivery of a certain QoS. To customers it serves as a warranty. An SLA usually include availability and performance guarantees. Additionally, metrics must be agreed upon by all parties as well as penalties for violating these expectations.

Most IaaS providers focus their SLA terms on availability guarantees, specifying the minimum percentage of time the system will be available during a certain period. For instance, Amazon EC2 states that —if the annual uptime Percentage for a customer drops below 99.95% for the service year, that customer is eligible to receive a service credit equal to 10% of their bill.³¹

Hypervisor and Operating System Choice— Traditionally, IaaS offerings have been based on heavily customized open-source Xen deployments. IaaS providers needed expertise in Linux, networking, virtualization, metering, resource management, and many other low-level aspects to successfully deploy and maintain their cloud offerings.

More recently, there has been an emergence of turnkey IaaS platforms such as VMWare VCloud and Citrix Cloud Center (C3) which have lowered the barrier of entry for IaaS competitors, leading to a rapid expansion in the IaaS marketplace.

Case Studies

Amazon Web Services: Amazon WS4 (AWS) is one of the major players in the cloud computing market. It pioneered the introduction of IaaS clouds in 2006. It offers a variety cloud services, most notably: S3 (storage), EC2 (virtual servers), Cloudfront (content delivery), Cloudfront Streaming (video streaming), Simple DB (structured datastore), RDS (Relational Database), SQS (reliable messaging), and Elastic MapReduce (data processing). The ElasticCompute Cloud (EC2) offers Xen-based virtual servers (instances) that can be

instantiated from Amazon Machine Images (AMIs). Instances are available in a variety of sizes, operating systems, architectures, and price. CPU capacity of instances is measured in Amazon Compute Units and, although fixed for each instance, vary among instance types from 1 (small instance) to 20 (high CPU instance). Each instance provides a certain amount of non persistent disk space; a persistence disk service (Elastic Block Storage) allows attaching virtual disks to instances with space up to 1TB. Elasticity can be achieved by combining the Cloud Watch, Auto Scaling and Elastic Load Balancing features, which allow the number of instances to scale up and down automatically based on a set of customizable rules, and traffic to be distributed across available instances. Fixed IP address (Elastic IPs) are not available by default, but can be obtained at an additional cost.

Flexiscale: Flexiscale is a UK-based provider offering services similar in nature to Amazon Web Services. Flexiscale cloud provides the following features: available in UK; Web services (SOAP), Web-based user interfaces; access to virtual server mainly via SSH (Linux) and Remote Desktop (Windows); 100% availability SLA with automatic recovery of VMs in case of hardware failure; per hour pricing; Linux and Windows operating systems; automatic scaling (horizontal/vertical).

Joyent: Joyent's Public Cloud offers servers based on Solaris containers virtualization technology. These servers, dubbed accelerators, allow deploying various specialized software- stack based on a customized version of Open- Solaris operating system, which include by default a Web-based configuration tool and several pre-installed software, such as Apache, MySQL, PHP, Ruby on Rails, and Java. Software load balancing is available as an accelerator in addition to hardware load balancers. A notable feature of Joyent's virtual servers is automatic vertical scaling of CPU cores, which means a virtual server can make use of additional CPUs automatically up to the maximum number of cores available in the physical host.

The Joyent public cloud offers the following features: multiple geographic locations in the United States; Web-based user interface; access to virtual server via SSH and Web-based administration tool; 100% availability SLA; per month pricing; OS-level virtualization Solaris containers; Open- Solaris operating systems; automatic scaling(vertical).

GoGrid: GoGrid, like many other IaaS providers, allows its customers to utilize a range of pre- made Windows and Linux images, in a range of fixed instance sizes. GoGrid also offers —value- added stacks on top for applications such as high- volume Web serving, e-Commerce, and database stores. It offers some notable features, such as a —hybrid hosting facility, which combines traditional dedicated hosts with auto-scaling cloud server infrastructure. As part of its core IaaS offerings, GoGrid also provides free hardware load balancing, auto-scaling capabilities, and persistent storage, features that typically add an additional cost for most other IaaS providers.

Rackspace Cloud Servers: Rackspace Cloud Servers is an IaaS solution that provides fixed size instances in the cloud. Cloud Servers offers a range of Linux- based pre-made images. A user can request different-sized images, where the size is measured by requested

RAM, not CPU.

PLATFORM AS A SERVICE PROVIDERS

Public Platform as a Service providers commonly offer a development and deployment environment that allow users to create and run their applications with little or no concern to low-level details of the platform. In addition, specific programming languages and frameworks are made available in the platform, as well as other services such as persistent data storage and in memory caches.

Features

Programming Models, Languages, and Frameworks: Programming models made available by IaaS providers define how users can express their applications using higher levels of abstraction and efficiently run them on the cloud platform.

Each model aims at efficiently solving a particular problem. In the cloud computing domain, the most common activities that require specialized models are: processing of large dataset in clusters of computers (MapReduce model), development of request-based Web services and applications; definition and orchestration of business processes in the form of workflows (Workflow model); and high-performance distributed execution of various computational tasks.

For user convenience, PaaS providers usually support multiple programming languages. Most commonly used languages in platforms include Python and Java (e.g., Google AppEngine), .NET languages (e.g., Microsoft Azure), and Ruby (e.g., Heroku). Force.com has devised its own programming language (Apex) and an Excel-like query language, which provide higher levels of abstraction to key platform functionalities.

A variety of software frameworks are usually made available to PaaS developers, depending on application focus. Providers that focus on Web and enterprise application hosting offer popular frameworks such as Ruby on Rails, Spring, Java EE, and .NET.

Persistence Options: A persistence layer is essential to allow applications to record their state and recover it in case of crashes, as well as to store user data. Web and enterprise application developers have chosen relational databases as the preferred persistence method. These databases offer fast and reliable structured data storage and transaction processing, but may lack scalability to handle several peta bytes of data stored in commodity computers. In the cloud computing domain, distributed storage technologies have emerged, which seek to be robust and highly scalable, at the expense of relational structure and convenient query languages.

CASE STUDIES

Aneka: Aneka is a .NET-based service-oriented resource management and development platform. Each server in an Aneka deployment (dubbed Aneka cloud node) hosts the Aneka container, which provides the base infrastructure that consists of services for persistence, security (authorization, authentication and auditing), and communication (message handling and dispatching). Cloud nodes can be either physical server, virtual machines (Xen Server and VMware are supported), and instances rented from Amazon EC2. The Aneka container can also host any number of optional services that can be added by developers to augment the capabilities of an Aneka Cloud node, thus providing a single, extensible framework for

orchestrating various application models.

Several programming models are supported by such task models to enable execution of legacy HPC applications and Map Reduce, which enables a variety of data-mining and search applications. Users request resources via a client to a reservation services manager of the Aneka master node, which manages all cloud nodes and contains scheduling service to distribute request to cloud nodes.

App Engine: Google App Engine lets you run your Python and Java Web applications on elastic infrastructure supplied by Google. App Engine allows your applications to scale dynamically as your traffic and data storage requirements increase or decrease. It gives developers a choice between a Python stack and Java. The App Engine serving architecture is notable in that it allows real-time auto- scaling without virtualization for many common types of Web applications. However, such auto-scaling is dependent on the application developer using a limited subset of the native APIs on each platform, and in some instances you need to use specific Google APIs such as URLFetch, Data store, and mem cache in place of certain native API calls. For example, a deployed App Engine application cannot write to the file system directly (you must use the Google Data store) or open a socket or access another host directly (you must use Google URL fetch service). A Java application cannot create a new Thread either.

Microsoft Azure: Microsoft Azure Cloud Services offers developers a hosted .NET Stack (C#, VB.Net, ASP.NET). In addition, a Java & Ruby SDK for .NET Services is also available. The Azure system consists of a number of elements. The Windows Azure Fabric Controller provides auto-scaling and reliability, and it manages memory resources and load balancing. The .NET Service Bus registers and connects applications together. The .NET Access Control identity providers include enterprise directories and Windows LiveID. Finally, the .NET Workflow allows construction and execution of workflow instances.

Force.com: In conjunction with the Salesforce.com service, the Force.com PaaS allows developers to create add-on functionality that integrates into main Salesforce CRM SaaS application. Force.com offers developers two approaches to create applications that can be deployed on its SaaS platform: a hosted Apex or Visualforce application. Apex is a proprietary Java-like language that can be used to create Salesforce applications. Visual force is an XML-like syntax for building UIs in HTML, AJAX, or Flex to overlay over the Salesforce hosted CRM system. An application store called App Exchange is also provided, which offers a paid & free application directory.

Heroku: Heroku is a platform for instant deployment of Ruby on Rails Web applications. In the Heroku system, servers are invisibly managed by the platform and are never exposed to users. Applications are automatically dispersed across different CPU cores and servers, maximizing performance and minimizing contention. Heroku has an advanced logic layer than can automatically route around failures, ensuring seamless and uninterrupted service at all times.

Public Cloud and Infrastructure Services

1. Public cloud or external cloud describes cloud computing in a traditional mainstream sense, whereby resources are dynamically provisioned via publicly accessible Web applications/Web services (SOAP or RESTful interfaces) from an off-site third-party provider.

2. Who shares resources and bills on a fine-grained utility computing basis, the user pays only for the capacity of the provisioned resources at a particular time.
3. Examples for vendors who publicly provide IaaS:
 1. Amazon Elastic Compute Cloud (EC2).
 2. GoGrid
 3. Joyent Accelerator
 4. Rackspace
 5. AppNexus
 6. FlexiScale and Manjrasoft Aneka
4. Amazon Elastic Compute Cloud (EC2) is an IaaS service that provides elastic compute capacity in the cloud.
5. These services can be leveraged via Web services (SOAP or REST), a Web-based AWS (Amazon Web Service) management console, or the EC2 command line tools.
6. The Amazon service provides hundreds of pre-made AMIs (Amazon Machine Images) with a variety of operating systems (i.e., Linux, OpenSolaris, or Windows) and pre-loaded software.
7. Provides complete control of computing resources run on Amazon's computing and infrastructure environment easily
8. Reduces the time required for obtaining and booting a new server's instances to minutes
9. Allows a quick scalable capacity and resources, up and down as the computing requirements change Offers different instances' size according to
 1. The resources' needs (small, large, and extra large)
 2. The high CPU's needs it provides (medium and extra large high CPU instances)
 3. High-memory instances (extra large, double extra large, and quadruple extra large instance)
10. Amazon EC2 is a widely known example for vendors that provide public cloud services.
11. Eucalyptus and Open-Nebula are two complementary and enabling technologies for open source cloud tools, which play an invaluable role in infrastructure as a service and in building private, public, and hybrid cloud architecture.
 1. The Amazon EC2 (Elastic Compute Cloud) is a Web service that allows users to provision new machines into Amazon's virtualized infrastructure in a matter of minutes using a publicly available API
 2. EC2 instance is typically a virtual machine with a certain amount of RAM, CPU, and storage capacity.

12. Amazon EC2 provides its customers with three flexible purchasing models to make it easy for the cost optimization:
 1. On-Demand instances: which allow you to pay a fixed rate by the hour with no commitment.
 2. Reserved instances: which allow you to pay a low, one-time fee and in turn receive a significant discount on the hourly usage charge for that instance. It ensures that any reserved instance you launch is guaranteed to succeed (provided that you have booked them in advance). This means that users of these instances should not be affected by any transient limitations in EC2 capacity.
 3. Spot instances: which enable you to bid whatever price you want for instance capacity, providing for even greater savings, if your applications have flexible start and end times.
13. Amazon Elastic Load Balancer is another service that helps in building fault-tolerant applications by automatically provisioning incoming application workload across available Amazon EC2 instances and in multiple availability zones.

Private Cloud and Infrastructure Services

A private cloud aims at providing public cloud functionality, but on private resources:

1. Maintaining control over an organization's data and resources to meet security and governance's requirements in an organization.
2. Private cloud exhibits a highly virtualized cloud data center located inside your organization's firewall.
3. It may also be a private space dedicated for your company within a cloud vendor's data center designed to handle the organization's workloads.

Private clouds exhibit the following characteristics:

1. Allow service provisioning and compute capability for an organization's users in a self-service manner.
2. Automate and provide well-managed virtualized environments.
3. Optimize computing resources, and servers' utilization.
4. Support specific workloads.

Examples for vendors and frameworks that provide IaaS in private setups

1. Eucalyptus (elastic utility computing architecture linking your programs to useful systems)
2. Open Nebula

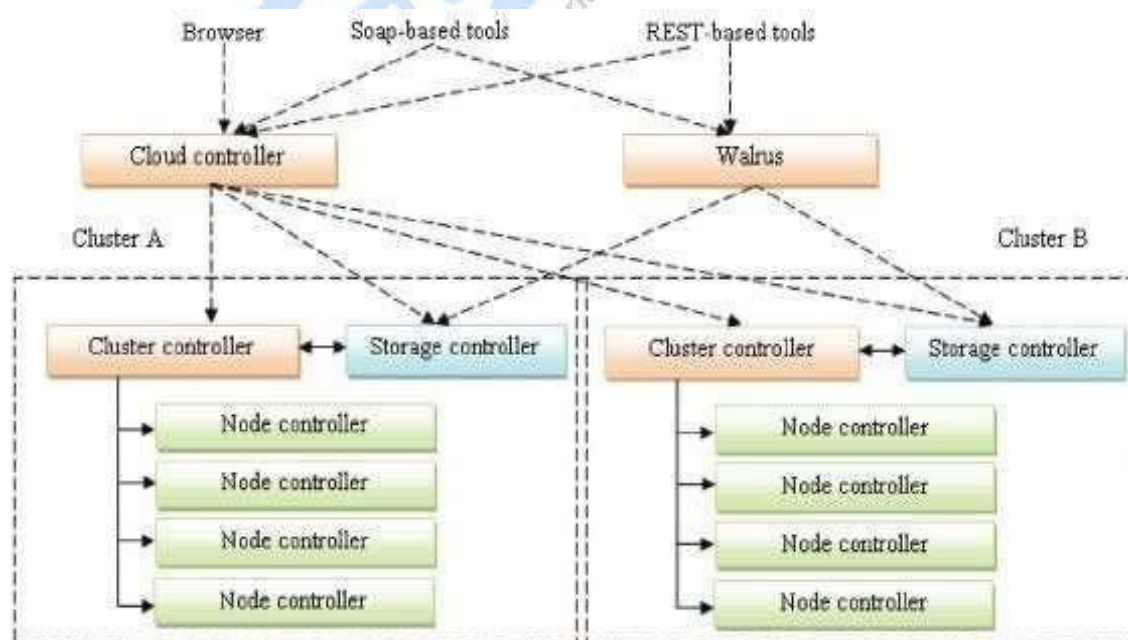
Eucalyptus: Eucalyptus is an open-source infrastructure for the implementation of cloud computing on computer clusters. It is considered one of the earliest tools developed as a surge

computing (in which data center's private cloud could augment its ability to handle workload's spikes by a design that allows it to send overflow work to a public cloud) tool. Its name is an acronym for —**elastic utility computing architecture for linking your programs to useful systems.**‖

Eucalyptus features :

1. Interface compatibility with EC2, and S3 (both Web service and Query/REST [Representational State Transfer] interfaces).
2. Simple installation and deployment.
3. Support for most Linux distributions (source and binary packages).
4. Support for running VMs that run atop the Xen hypervisor or KVM.
5. Support for other kinds of VMs, such as VMware, is targeted for future releases.
6. Secure internal communication using SOAP (Simple Object Access Protocol) with WS security.
7. Cloud administrator's tool for system's management and user's accounting.
8. The ability to configure multiple clusters each with private internal network addresses into a single cloud.
9. Eucalyptus aims at fostering the research in models for service's provisioning, scheduling, SLA formulation, and hypervisors' portability.

Eucalyptus Architecture:



1. **Node controller (NC)** controls the execution, inspection, and termination of VM instances on the host where it runs.
2. **Cluster controller (CC)** gathers information about and schedules VM execution on specific node controllers, as well as manages virtual instance network.

3. **Storage controller (SC)** is a put/get storage service that implements Amazon's S3(Simple Storage Service) interface and provides a way for storing and accessing VM images and user data.
4. **Cloud controller (CLC)** is the entry point into the cloud for users and administrators. It queries node managers for information about resources, makes high-level scheduling decisions, and implements them by making requests to cluster controllers.
5. **Walrus (W)** is the controller component that manages access to the storage services within Eucalyptus. Requests are been communicated to Walrus using the SOAP (Simple Object Access Protocol) or REST (Representational State Transfer) based interface

Hybrid Cloud and Infrastructure Services

A third type of cloud setup named Hybrid cloud

1. A combination of private/internal and external cloud resources existing together by enabling outsourcing of noncritical services and functions in public cloud and keeping the critical ones internal
2. Main function of Hybrid cloud is to release resources from a public cloud and handle sudden demand usage called cloud bursting

Distributed Management of Virtualization

Virtualization's benefits bring their own challenges and complexities presented in the need for a powerful management capabilities. That is why many commercial, open source products and research projects such as OpenNebula, IBM Virtualization Manager, Joyent, and VMware DRS are been developed to be dynamically provision virtual machines, utilizing the physical infrastructure. There are also some commercial and scientific infrastructure cloud computing initiatives, such as Globus VWS, Eucalyptus and Amazon, which provide remote interfaces for controlling and monitoring virtual resources.

One more effort in this context is the RESERVOIR initiative, in which grid interfaces and protocols enable the required interoperability between the clouds or infrastructure's providers.

High Availability

High availability is a system design protocol and an associated implementation that ensures a certain absolute degree of operational continuity during a given measurement period. Availability refers to the ability of a user's community to access the system—whether for submitting new work, updating or altering existing work, or collecting the results of the previous work.

Cloud and Virtualization Standardization Efforts

Standardization is important to ensure interoperability between virtualization management vendors, the virtual machines produced by each one of them, and cloud computing. In the past few years, virtualization standardization efforts led by the Distributed

Management Task Force (DMTF) have produced standards for almost all the aspects of virtualization technology.

DMTF initiated the VMAN (Virtualization Management Initiative), which delivers broadly supported interoperability and portability standards for managing the virtual computing lifecycle. VMAN's OVF (Open Virtualization Format) is a collaboration between industry key players: Dell, HP, IBM, Microsoft, XenSource, and VMware.

OVF (Open Virtualization Format)

1. VMAN's OVF (Open Virtualization Format) is a collaboration between industry key players: Dell, HP, IBM, Microsoft, XenSource, and VMware.
2. OVF specification provides a common format to package and securely distribute virtual appliances across multiple virtualization platforms.
3. VMAN profiles define a consistent way of managing a heterogeneous virtualized environment

OCCI and OGF

Open Grid Forum (OGF) organizing an official new working group to deliver a standard API for cloud IaaS, the Open Cloud Computing Interface Working Group (OCCIWG). This group is dedicated for delivering an API specification for the remote management of cloud computing's infrastructure and for allowing the development of interoperable tools for common tasks including deployment, autonomic scaling, and monitoring. The scope of the specification will be covering a high-level functionality required for managing the life-cycle virtual machines (or workloads), running on virtualization technologies (or containers), and supporting service elasticity. The new API for interfacing —IaaS cloud computing facilities will allow

1. **Consumers** to interact with cloud computing infrastructure on an ad hoc basis.
2. **Integrators** to offer advanced management services.
3. **Aggregators** to offer a single common interface to multiple providers. Providers to offer a standard interface that is compatible with the available tools.
4. **Vendors** of grids/clouds to offer standard interfaces for dynamically scalable service's delivery in their products.

VM Provisioning Process

Typical life cycle of VM and its major possible states of operation, which make the management and automation of VMs in virtual and cloud environments easier. Process & Steps to Provision VM. Here, we describe the common and normal steps of provisioning a virtual server:

1. Firstly, you need to select a server from a pool of available servers (physical servers with enough capacity) along with the appropriate OS template you need to provision the virtual machine.
2. Secondly, you need to load the appropriate software (operating system you selected in the previous step, device drivers, middleware, and the needed applications for the service required).

3. Thirdly, you need to customize and configure the machine (e.g., IP address, Gateway) to configure an associated network and storage resources.
4. Finally, the virtual server is ready to start with its newly loaded software. Typically, these are the tasks required or being performed by an IT or a data center's specialist to provision a particular virtual machine.

virtual machines can be provisioned by manually installing an operating system, by using a preconfigured VM template, by cloning an existing VM, or by importing a physical server or a virtual server from another hosting platform. Physical servers can also be virtualized and provisioned using P2V (physical to virtual) tools and techniques (e.g., virt-p2v).

After creating a virtual machine by virtualizing a physical server, or by building a new virtual server in the virtual environment, a template can be created out of it. Most virtualization management vendors (VMware, XenServer, etc.) provide the data center's administration with the ability to do such tasks in an easy way.

Provisioning from a template is an invaluable feature, because it reduces the time required to create a new virtual machine. Administrators can create different templates for different purposes. For example, you can create a Windows 2003 Server template for the finance department, or a Red Hat Linux template for the engineering department.

This enables the administrator to quickly provision a correctly configured virtual server on demand. This ease and flexibility bring with them the problem of virtual machine's sprawl, where virtual machines are provisioned so rapidly that documenting and managing

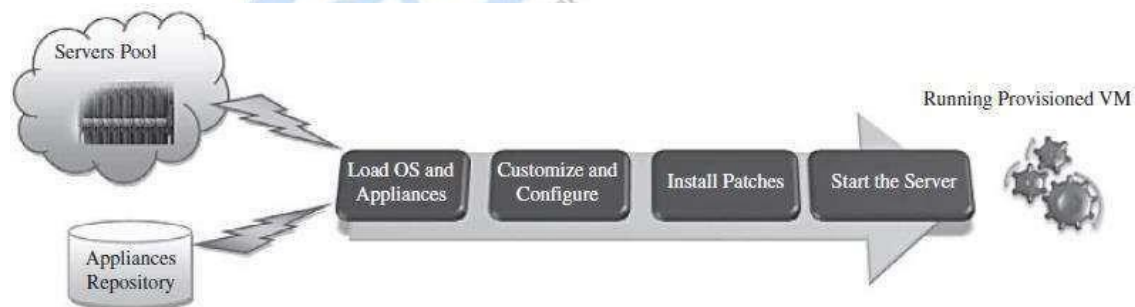


FIGURE 5.4. Virtual machine provision process.

the virtual machine's life cycle become a challenge

VIRTUAL MACHINE MIGRATION SERVICES

Migration service, in the context of virtual machines, is the process of moving a virtual machine from one host server or storage location to another; there are different techniques of VM migration, hot/live migration, cold/regular migration, and live storage migration of a virtual machine. In this process, all key machine components, such as CPU, storage disks, networking, and memory, are completely virtualized, thereby facilitating the entire state of a virtual machine to be captured by a set of easily moved data files. Here are some of the migration's techniques that most virtualization tools provide as a feature.

Migrations Techniques

Live Migration and High Availability.

Live migration (which is also called hot or real-time migration) can be defined as the movement of a virtual machine from one physical host to another while being powered on. When it is properly carried out, this process takes place without any noticeable effect from the end user's point of view (a matter of milliseconds). One of the most significant advantages of live migration is the fact that it facilitates proactive maintenance in case of failure, because the potential problem can be resolved before the disruption of service occurs. Live migration can also be used for load balancing in which work is shared among computers in order to optimize the utilization of available CPU resources.

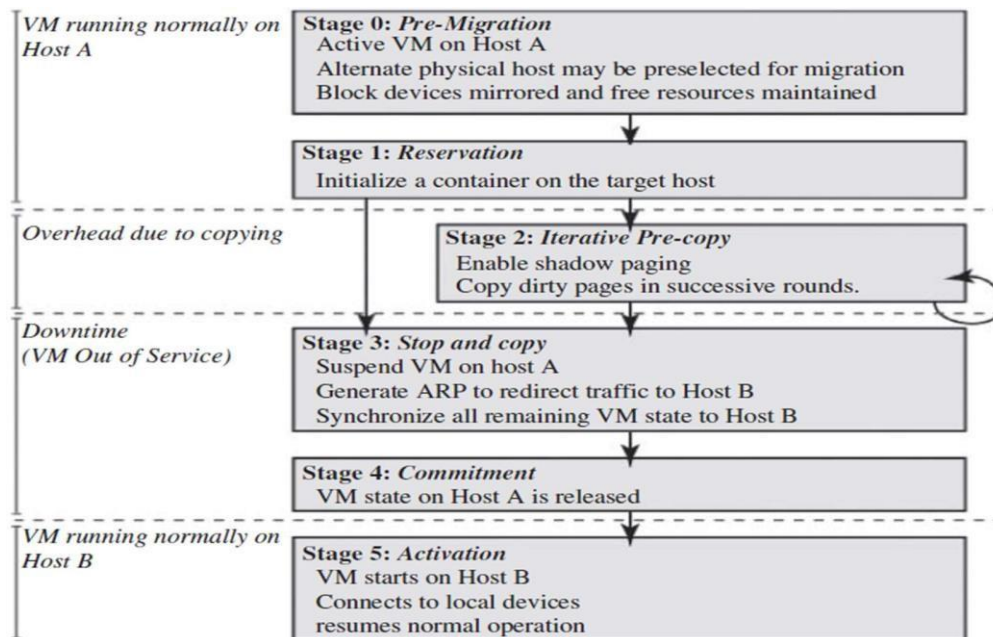
Live Migration Anatomy, Xen Hypervisor Algorithm.

In this section we will explain live migration's mechanism and how memory and virtual machine states are being transferred, through the network, from one host A to another host B, the Xen hypervisor is an example for this mechanism. The logical steps that are executed when migrating an OS are summarized in the diagram below. In this research, the migration process has been viewed as a transactional interaction between the two hosts involved.

Migration Techniques:

- Stage 0: Pre-Migration
 - An active virtual machine exists on the physical host A
- Stage 1: Reservation
 - A request is issued to migrate an OS from host A to B.
 - The necessary resources exist on B and on a VM container of that size.
- Stage 2: Iterative Pre-Copy
 - During the first iteration, all pages are transferred from A to B
 - Subsequent iterations copy only those pages dirtied during the previous transfer phase
- Stage 3: Stop-and-Copy
 - Running OS instance at A is suspended
 - The network traffic is redirected to B
 - CPU state and any remaining inconsistent memory pages are then transferred
 - At the end of this stage, there is a consistent suspended copy of the VM at both A and B.
 - Copy at A is considered primary and is resumed in case of failure
- Stage 4: Commitment
 - Host B indicates to A that it has successfully received a consistent OS image
 - Host A acknowledges this message as a commitment of the migration transaction
 - Host A may now discard the original VM
 - Host B becomes the primary host
- Stage 5: Activation
 - The migrated VM on B is now activated

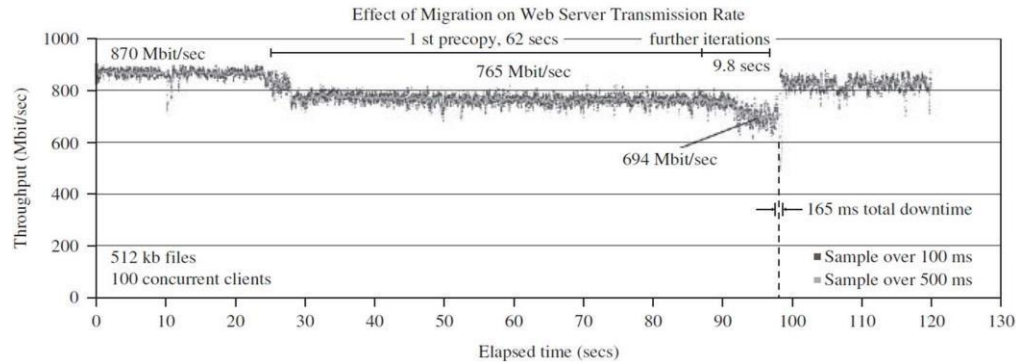
- Post-migration code runs to reattach the device's drivers to the new machine and advertise moved IP addresses



Live Migration Effect on a Running Web Server.

Clark et al. Did evaluate the above migration on an Apache 1.3 Web server; this served static content at a high rate as shown in diagram below. The throughput is achieved when continuously serving a single 512-kB file to a set of one hundred concurrent clients. The Web server virtual machine has a memory allocation of 800 MB. At the start of the trace, the server achieves a consistent throughput of approximately 870 Mbit/sec.

Migration starts 27 sec into the trace, but is initially rate-limited to 100 Mbit/sec (12% CPU), resulting in server's throughput drop to 765 Mbit/sec. This initial low-rate pass transfers 776 MB and lasts for 62 sec. At this point, the migration's algorithm, described in Section 5.4.1, increases its rate over several iterations and finally suspends the VM after a further 9.8 sec. The final stop-and-copy phase then transfers the remaining pages, and the Web server resumes at full rate after a 165-msec outage. This simple example demonstrates that a highly loaded server can be migrated with both controlled impact on live services and a short downtime. However, the working set of the server, in this case, is rather small. So, this should be expected as a relatively easy case of live migration.



Live Migration Vendor Implementations Examples.

There are lots of VM management and provisioning tools that provide the live migration of VM facility, two of which are VMware VMotion and Citrix XenServer —XenMotion.¶

VMware Vmotion.

This allows users to automatically optimize and allocate an entire pool of resources for maximum hardware utilization, flexibility, and availability and perform hardware's maintenance without scheduled downtime along with migrating virtual machines away from failing or underperforming servers.

Citrix XenServer XenMotion.

This is a nice feature of the Citrix XenServer product, inherited from the Xen live migrate utility, which provides the IT administrator with the facility to move a running VM from one XenServer to another in the same pool without interrupting the service (hypothetically for zero-downtime server maintenance, which actually takes minutes), making it a highly available service. This also can be a good feature to balance the workloads on the virtualized environment.

Regular/Cold Migration.

Cold migration is the migration of a powered-off virtual machine. With cold migration, you have the option of moving the associated disks from one data store to another. The virtual machines are not required to be on a shared storage.

It's important to highlight that the two main differences between live migration and cold migration are that live migration needs a shared storage for virtual machines in the server's pool, but cold migration does not and also in live migration for a virtual machine between two hosts, there would be certain CPU compatibility checks to be applied; while in cold migration this checks do not apply.

The cold migration process is simple to implement and it can be summarized as follows:

1. The configuration files, including the NVRAM file (BIOS settings), log files, as well as the disks of the virtual machine, are moved from the source host to the destination host's associated storage area.
2. The virtual machine is registered with the new host.
3. After the migration is completed, the old version of the virtual machine is deleted from the source host.

Live Storage Migration of Virtual Machine.

This kind of migration constitutes moving the virtual disks or configuration file of a running virtual machine to a new data store without any interruption in the availability of the virtual machine's service.

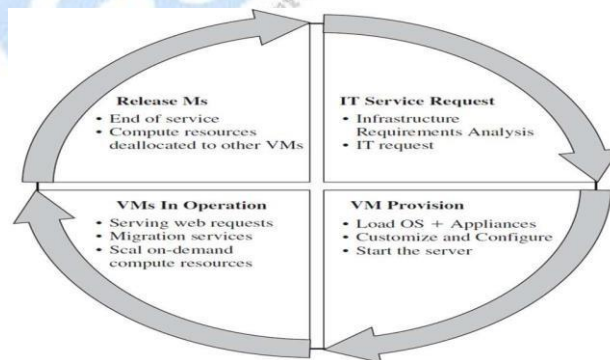
Migration of Virtual Machines to Alternate Platforms

One of the nicest advantages of having facility in data center's technologies is to have the ability to migrate virtual machines from one platform to another. There are a number of ways for achieving this, such as depending on the source and target virtualization's platforms and on the vendor's tools that manage this facility—for example, the VMware converter that handles migrations between ESX hosts; the VMware server; and the VMware workstation. The VMware converter can also import from other virtualization platforms, such as Microsoft virtual server machines.

VIRTUAL MACHINES PROVISIONING AND MANAGEABILITY

The typical life cycle of VM and its major possible states of operation, which make the management and automation of VMs in virtual and cloud environments easier than in traditional computing environments.

As shown in the diagram below the cycle starts by a request delivered to the IT department, stating the requirement for creating a new server for a particular service. This request is being processed by the IT administration to start seeing the servers' resource pool, matching these resources with the requirements, and starting the provision of the needed virtual machine. Once it is provisioned and started, it is ready to provide the required service according to an SLA, or a time period after which the virtual is being released; and free resources, in this case, won't be needed.



VM PROVISIONING AND MIGRATION IN ACTION

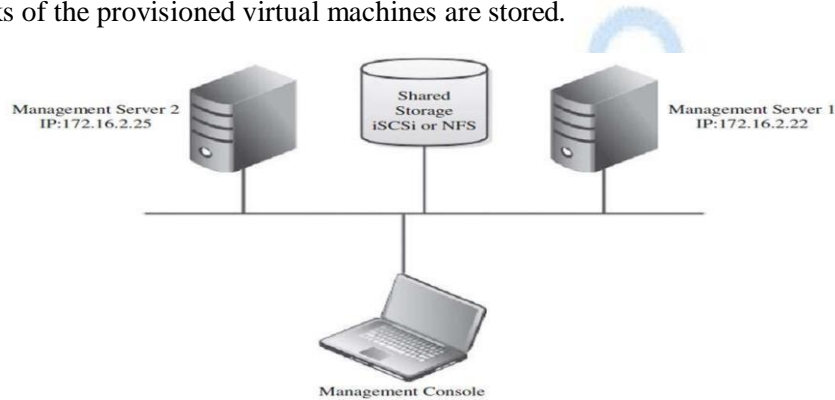
Now, it is time to get into business with a real example of how we can manage the life cycle, provision, and migrate a virtual machine by the help of one of the open source frameworks used to manage virtualized infrastructure.

Here, we will use ConVirt (open source framework for the management of open source virtualization like Xen and KVM known previously as XenMan). Deployment Scenario. ConVirt deployment consists of at least one ConVirt workstation, where ConVirt is installed and ran, which provides the main console for managing the VM life cycle, managing images, provisioning new VMs, monitoring machine resources, and so on.

There are two essential deployment scenarios for ConVirt:

1. Basic configuration in which the Xen or KVM virtualization platform is on the local machine, where ConVirt is already installed.
2. An advanced configuration in which the Xen or KVM is on one or more remote servers. The scenario in use here is the advanced one. In data centers, it is very common to install centralized management software (ConVirt here) on a dedicated machine for use in managing remote servers in the data center.

In our example, we will use this dedicated machine where ConVirt is installed and used to manage a pool of remote servers (two machines). In order to use advanced features of ConVirt (e.g., live migration), you should set up a shared storage for the server pool in use on which the disks of the provisioned virtual machines are stored.



Installation

The installation process involves the following:

1. Installing ConVirt on at least one computer. See reference 28 for installation details.
2. Preparing each managed server to be managed by ConVirt. See reference 28 for managed servers' installation details. We have two managing servers with the following Ips (managed server 1, IP:172.16.2.22; and managed server 2, IP:172.16.2.25) as shown in the deployment diagram (Figure 5.7).
3. Starting ConVirt and discovering the managed servers you have prepared.

Notes

1. Try to follow the installation steps existing in reference 28 according to the distribution of the operating system in use. In our experiment, we use Ubuntu 8.10 in our setup.
2. Make sure that the managed servers include Xen or KVM hypervisors installed.
3. Make sure that you can access managed servers from your ConVirt management console through SSH.

Environment, Software, and Hardware: ConVirt 1.1, Linux Ubuntu 8.10, three machines, Dell core 2 due processor, 4G RAM.

Adding Managed Servers and Provisioning VM: Once the installation is done and you are ready to manage your virtual infrastructure, then you can start the ConVirt management console.

Select any of servers' pools existing (QA Lab in our scenario) and on its context menu, select —Add Serverl.

1. You will be faced with a message asking about the virtualization platform you want to manage (Xen or KVM).
2. Choose KVM, and then enter the managed server information and credentials (IP, username, and password).
3. Once the server is synchronized and authenticated with the management console, it will appear in the left pane/of the ConVirt.
4. Select this server, and start provisioning your virtual machine.
5. Fill in the virtual machine's information (name, storage, OS template, etc) then you will find it created on the managed server tree powered-off. Note: While provisioning your virtual machine, make sure that you create disks on the shared storage (NFS or iSCSi). You can do so by selecting the —provisioningl tab, and changing the VM_DISKS_DIR to point to the location of your shared NFS.
6. Start your VM (Figures 5.14 and 5.15), and make sure the installation media of the operating system you need is placed in drive, in order to use it for booting the new VM and proceed in the installation process; then start the installation process as shown in Figure 5.16.
7. Once the installation finishes, you can access your provisioned virtual machine from the console icon on the top of your ConVirt management console.
8. Reaching this step, you have created your first managed server and provisioned virtual machine. You can repeat the same procedure to add the second managed server in your pool to be ready for the next step of migrating one virtual machine from one server to the other.
9. To start the migration of a virtual machine from one host to the other, select it and choose a migrating virtual machine.
10. You will have a window containing all the managed servers in your data center. Choose one as a destination and start
11. Once the virtual machine has been successfully placed and migrated to the destination host, you can see it still living and working.

ON THE MANAGEMENT OF VIRTUAL MACHINES FOR CLOUD INFRASTRUCTURES

In 2006, Amazon started offering virtual machines (VMs) to anyone with a credit card for just \$0.10/hour through its Elastic Compute Cloud (EC2) service. Although not the first company to lease VMs, the programmer-friendly EC2 Web services API and their pay-as-you-go pricing popularized the —Infrastructure as a Servicel (IaaS) paradigm, which is now closely related to the notion of a —cloud.l

Following the success of Amazon EC2, several other IaaS cloud providers, or public clouds, have emerged—such as Elastic- Hosts, GoGrid, and FlexiScale—that provide a publicly

accessible interface for purchasing and managing computing infrastructure that is instantiated as VMs running on the provider's data center.

There is also a growing ecosystem of technologies and tools to build private clouds—where inhouse resources are virtualized, and internal users can request and manage these resources using interfaces similar or equal to those of public clouds—and hybrid clouds—where an organization's private cloud can supplement its capacity using a public cloud.

THE ANATOMY OF CLOUD INFRASTRUCTURES

There are many commercial IaaS cloud providers in the market, such as those cited earlier, and all of them share five characteristics:

- (i) They provide on-demand provisioning of computational resources.
 - (ii) they use virtualization technologies to lease these resources.
 - (iii) they provide public and simple remote interfaces to manage those resources
 - (iv) they use a pay-as-you-go cost model, typically charging by the hour
 - (v) they operate data centers large enough to provide a seemingly unlimited amount of resources to their clients (usually touted as —infinite capacity‖ or —unlimited elasticity‖).
1. Private and hybrid clouds share these same characteristics but, instead of selling capacity over publicly accessible interfaces, focus on providing capacity to an organization's internal users.
 2. Virtualization technologies have been the key enabler of many of these salient characteristics of IaaS clouds by giving providers a more flexible and generic way of managing their resources. Thus, virtual infrastructure (VI) management—the management of virtual machines distributed across a pool of physical resources—becomes a key concern when building an IaaS cloud and poses a number of challenges.
 3. Virtual infrastructure management in private clouds has to deal with an additional problem: Unlike large IaaS cloud providers, such as Amazon, private clouds typically do not have enough resources to provide the illusion of —infinite capacity.‡ The immediate provisioning scheme used in public clouds, where resources are provisioned at the moment they are requested, is ineffective in private clouds.
 4. Several VI management solutions have emerged over time, such as platform ISF and VMware vSphere, along with open-source initiatives such as Enomaly Computing Platform and Ovirt.
 5. However, managing virtual infrastructures in a private/hybrid cloud is a different, albeit similar, problem than managing a virtualized data center, and existing tools lack several features that are required for building IaaS clouds.

Distributed Management of Virtual Machines

The first problem is how to manage the virtual infrastructures themselves. Although resource management has been extensively studied, particularly for job management in high-performance computing, managing VMs poses additional problems that do not arise when managing jobs, such as the need to set up custom software environments for VMs, setting up and managing networking for interrelated VMs, and reducing the various overheads involved in using VMs.

1. Thus, VI managers must be able to efficiently orchestrate all these different tasks. The problem of efficiently selecting or scheduling computational resources is well known.
2. However, the state of the art in VM-based resource scheduling follows a static approach, where resources are initially selected using a greedy allocation strategy, with minimal or no support for other placement policies.
3. To efficiently schedule resources, VI managers must be able to support flexible and complex scheduling policies and must leverage the ability of VMs to suspend, resume, and migrate. This complex task is one of the core problems that the RESERVOIR (Resources and Services Virtualization without Barriers) project tries to solve.

Reservation-Based Provisioning of Virtualized Resources

A particularly interesting problem when provisioning virtual infrastructures is how to deal with situations where the demand for resources is known beforehand— for example, when an experiment depending on some complex piece of equipment is going to run from 2 pm to 4 pm, and computational resources must be available at exactly that time to process the data produced by the equipment. Commercial cloud providers, such as Amazon, have enough resources to provide the illusion of infinite capacity, which means that this situation is simply resolved by requesting the resources exactly when needed; if capacity is —infinite, then there will be resources available at 2 pm. On the other hand, when dealing with finite capacity, a different approach is needed. However, the intuitively simple solution of reserving the resources beforehand turns out to not be so simple, because it is known to cause resources to be underutilized, due to the difficulty of scheduling other requests around an inflexible reservation. VMs allow us to overcome the utilization problems typically associated with advance reservations and we describe Haizea, a VM-based lease manager supporting advance reservation along with other provisioning models not supported in existing IaaS clouds, such as best-effort provisioning.

Provisioning to Meet SLA Commitments

IaaS clouds can be used to deploy services that will be consumed by users other than the one that deployed the services. For example, a company might depend on an IaaS cloud provider to deploy three-tier applications (Web front-end, application server, and database server) for its customers. In this case, there is a distinction between the cloud consumer (i.e., the service owner) and the end users of the resources provisioned on the cloud (the service user).

Furthermore, service owners will enter into service-level agreements (SLAs) with their end users, covering guarantees such as the timeliness with which these services will

respond. However, cloud providers are typically not directly exposed to the service semantics or the SLAs that service owners may contract with their end users. The capacity requirements are less predictable and more elastic.

The cloud provider's task is, therefore, to make sure that resource allocation requests are satisfied with specific probability and timeliness. These requirements are formalized in infrastructure SLAs between the service owner and cloud provider, separate from the high-level SLAs between the service owner and its end users.

RESERVOIR proposes a flexible framework where service owners may register service-specific elasticity rules and monitoring probes, and these rules are being executed to match environment conditions.

Elasticity of the application should be contracted and formalized as part of capacity availability SLA between the cloud provider and service owner. This poses interesting research issues on the IaaS side, which can be grouped around two main topics:

1. SLA-oriented capacity planning that guarantees that there is enough capacity to guarantee service elasticity with minimal over-provisioning.
2. Continuous resource placement and scheduling optimization that lowers operational costs and takes advantage of available capacity transparently to the service while keeping the service SLAs.

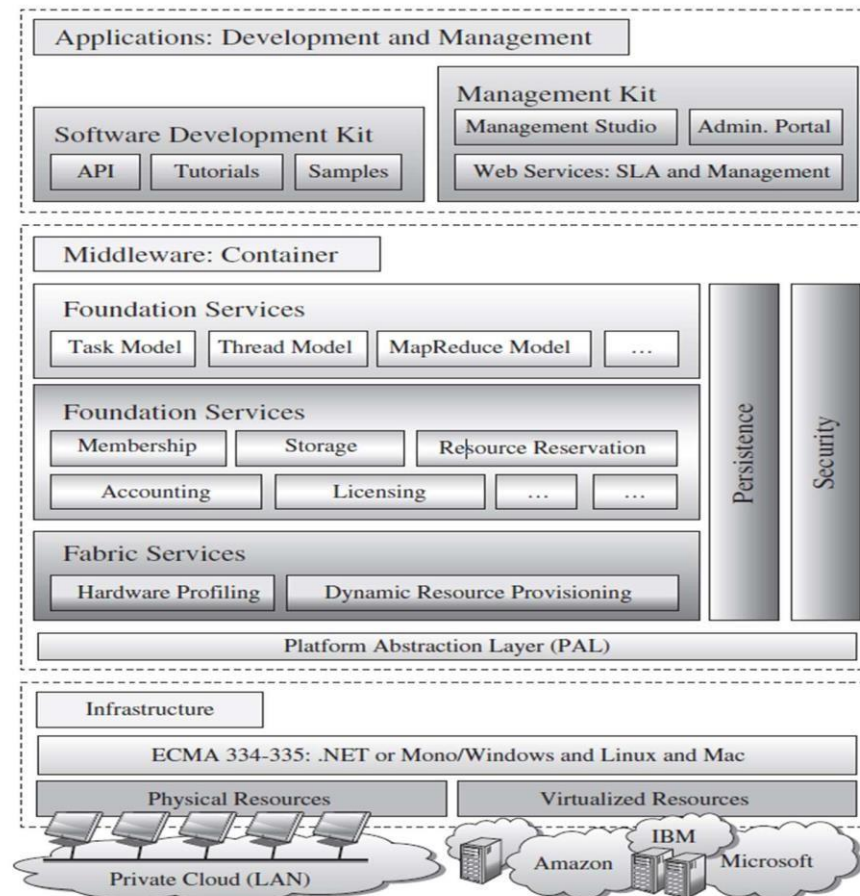
ANEKA—INTEGRATION OF PRIVATE AND PUBLIC CLOUDS

1. Aneka is a software platform and a framework for developing distributed applications on the cloud. It harnesses the computing resources of a heterogeneous network of workstations and servers or data centers on demand. Aneka provides developers with a rich set of APIs for transparently exploiting these resources by expressing the application logic with a variety of programming abstractions. System administrators can leverage a collection of tools to monitor and control the deployed infrastructure.
2. This can be a public cloud available to anyone through the Internet, a private cloud constituted by a set of nodes with restricted access within an enterprise, or a hybrid cloud where external resources are integrated on demand, thus allowing applications to scale. Diagram below provides a layered view of the framework.
3. Aneka is essentially an implementation of the PaaS model, and it provides a runtime environment for executing applications by leveraging the underlying infrastructure of the cloud. Developers can express distributed applications by using the API contained in the Software Development Kit (SDK) or by porting existing legacy applications to the cloud.
4. Such applications are executed on the Aneka cloud, represented by a collection of nodes connected through the network hosting the Aneka container.
5. The container is the building block of the middleware and represents the runtime environment for executing applications; it contains the core functionalities of the

system and is built up from an extensible collection of services that allow administrators to customize the Aneka cloud. There are three classes of services that characterize the container:

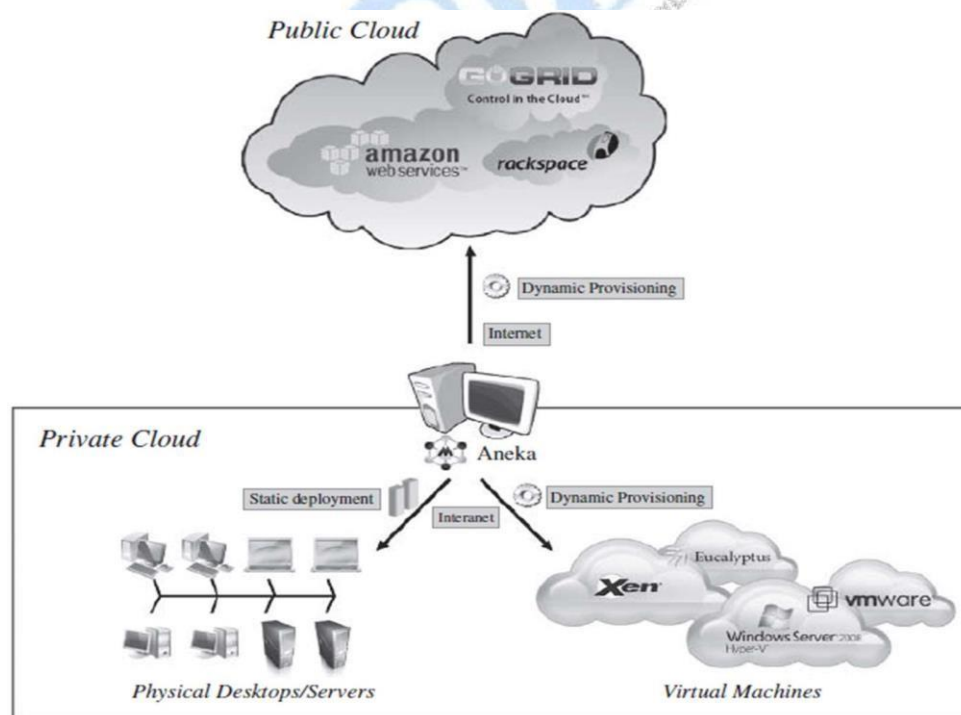
6. **Execution Services.** They are responsible for scheduling and executing applications. Each of the programming models supported by Aneka defines specialized implementations of these services for managing the execution of a unit of work defined in the model.
7. **Foundation Services.** These are the core management services of the Aneka container. They are in charge of metering applications, allocating resources for execution, managing the collection of available nodes, and keeping the services registry updated.
8. **Fabric Services:** They constitute the lowest level of the services stack of Aneka and provide access to the resources managed by the cloud. An important service in this layer is the Resource Provisioning Service, which enables horizontal scaling³ in the cloud. Resource provisioning makes Aneka elastic and allows it to grow or to shrink dynamically to meet the QoS requirements of applications.
9. Aneka also provides a tool for managing the cloud, allowing administrators to easily start, stop, and deploy instances of the Aneka container on new resources and then reconfigure them dynamically to alter the behavior of the cloud.

Aneka Resource Provisioning Service
1.



uting has the ability to automatically scale out

1. Based on demand and users' quality of service requests
2. Aneka is a PaaS
 1. Features multiple programming models allowing developers to easily build their distributed applications
 2. Provides resource provisioning facilities in a seamless and dynamic fashion
 3. Achieved by means of the resource provisioning framework
3. A typical scenario that a medium or large enterprise may encounter
 1. Combines privately owned resources with public rented resources to dynamically increase the resource capacity to a larger scale
4. Private resources identify computing and storage elements kept in the premises. They share similar internal security and administrative policies



CHALLENGES AND RISKS:

Despite the initial success and popularity of the cloud computing paradigm and the extensive availability of providers and tools, a significant number of challenges and risks are inherent to this new model of computing. Providers, developers, and end users must consider these challenges and risks to take good advantage of cloud computing. Issues to be faced include

user privacy, data security, data lock-in, availability of service, disaster recovery, performance, scalability, energy-efficiency, and programmability.

Security, Privacy, and Trust: Security and privacy affect the entire cloud computing stack, since there is a massive use of third-party services and infrastructures that are used to host important data or to perform critical operations. In this scenario, the trust toward providers is fundamental to ensure the desired level of privacy for applications hosted in the cloud. Legal and regulatory issues also need attention. When data are moved into the Cloud, providers may choose to locate them anywhere on the planet. The physical location of data centers determines the set of laws that can be applied to the management of data. For example, specific cryptography techniques could not be used because they are not allowed in some countries. Similarly, country laws can impose that sensitive data, such as patient health records, are to be stored within national borders.

Data Lock-In and Standardization: A major concern of cloud computing users is about having their data locked-in by a certain provider. Users may want to move data and applications out from a provider that does not meet their requirements. However, in their current form, cloud computing infrastructures and platforms do not employ standard methods of storing user data and applications. Consequently, they do not interoperate and user data are not portable.

The answer to this concern is standardization. In this direction, there are efforts to create open standards for cloud computing. The Cloud Computing Interoperability Forum (CCIF) was formed by organizations such as Intel, Sun, and Cisco in order to —enable a global cloud computing ecosystem whereby organizations are able to seamlessly work together for the purposes for wider industry adoption of cloud computing technology. The development of the Unified Cloud Interface (UCI) by CCIF aims at creating a standard programmatic point of access to an entire cloud infrastructure. In the hardware virtualization sphere, the Open Virtual Format (OVF) aims at facilitating packing and distribution of software to be run on VMs so that virtual appliances can be made portable—that is, seamlessly run on hypervisor of different vendors.

Availability, Fault-Tolerance, and Disaster Recovery: It is expected that users will have certain expectations about the service level to be provided once their applications are moved to the cloud. These expectations include availability of the service, its overall performance, and what measures are to be taken when something goes wrong in the system or its components. In summary, users seek for a warranty before they can comfortably move their business to the cloud. SLAs, which include QoS requirements, must be ideally set up between customers and cloud computing providers to act as warranty. An SLA specifies the details of the service to be provided, including availability and performance guarantees. Additionally, metrics must be agreed upon by all parties, and penalties for violating the expectations must also be approved.

Resource Management and Energy-Efficiency: One important challenge faced by providers of cloud computing services is the efficient management of virtualized resource pools. Physical resources such as CPU cores, disk space, and network bandwidth must be sliced and shared among virtual machines running potentially heterogeneous workloads. The

multi-dimensional nature of virtual machines complicates the activity of finding a good mapping of VMs onto available physical hosts while maximizing user utility. Dimensions to be considered include: number of CPUs, amount of memory, size of virtual disks, and network bandwidth.

Dynamic VM mapping policies may leverage the ability to suspend, migrate, and resume VMs as an easy way of preempting low-priority allocations in favor of higher-priority ones. Migration of VMs also brings additional challenges such as detecting when to initiate a migration, which VM to migrate, and where to migrate. In addition, policies may take advantage of live migration of virtual machines to relocate data center load without significantly disrupting running services.

In this case, an additional concern is the trade-off between the negative impact of a live migration on the performance and stability of a service and the benefits to be achieved with that migration. Another challenge concerns the outstanding amount of data to be managed in various VM management activities. Such data amount is a result of particular abilities of virtual machines, including the ability of traveling through space (i.e., migration) and time (i.e., check pointing and rewinding), operations that may be required in load balancing, backup, and recovery scenarios. In addition, dynamic provisioning of new VMs and replicating existing VMs require efficient mechanisms to make VM block storage devices (e.g., image files) quickly available at selected hosts. Data centers consumer large amounts of electricity. According to a data published by HP, 100 server racks can consume 1.3MW of power and another 1.3 MW are required by the cooling system, thus costing USD 2.6 million per year. Besides the monetary cost, data centers significantly impact the environment in terms of CO₂ emissions from the cooling systems

UNIT – IV




Software as a Service (SAAS) & Data Security in the Cloud: Software as a Service SAAS), Google App Engine – Centralizing Email Communications- Collaborating via Web- Based Communication Tools-An Introduction to the idea of Data Security.

Software as a Service (SAAS) & Data Security in the Cloud

Software as a Service(SAAS)

Traditional desktop applications such as word processing and spreadsheet can now be accessed as a service in the Web. This model of delivering applications, known as Software as a Service (SaaS), alleviates the burden of software maintenance for customers and simplifies development and testing for providers.

Salesforce.com, which relies on the SaaS model, offers business productivity applications (CRM) that reside completely on their servers, allowing customers to customize and access applications on demand.

Service Class	Main Access & Management Tool	Service content
 SaaS	Web Browser	Cloud Applications Social networks, Office suites, CRM, Video processing
 PaaS	Cloud Development Environment	Cloud Platform Programming languages, Frameworks, Mashups editors, Structured data
 IaaS	Virtual Infrastructure Manager	Cloud Infrastructure Compute Servers, Data Storage, Firewall, Load Balancer

17

Deployment Models

Although cloud computing has emerged mainly from the appearance of public computing utilities, other deployment models, with variations in physical location and distribution, have been adopted. In this sense, regardless of its service class, a cloud can be classified as public, private, community, or hybrid based on model of

deployment as shown figure below.

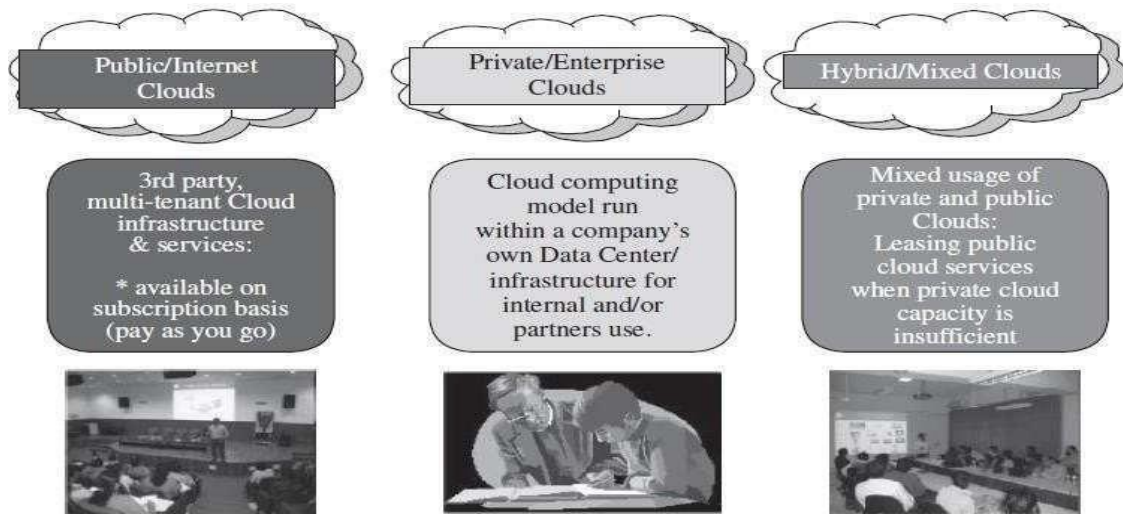


FIGURE 1.4. Types of clouds based on deployment models.

Public cloud & Private cloud:

Public cloud as a —cloud made available in a pay-as-you-go manner to the general public. **Private cloud** as —internal data center of a business or other organization, not made available to the general public.

In most cases, establishing a private cloud means restructuring an existing infrastructure by adding virtualization and cloud-like interfaces. This allows users to interact with the local data center while experiencing the same advantages of public clouds, most notably self-service interface, privileged access to virtual servers, and per-usage metering and billing.

A **community cloud** is —shared by several organizations and supports a specific community that has shared concerns (e.g., mission, security requirements, policy, and compliance considerations).

A **hybrid cloud** takes shape when a private cloud is supplemented with computing capacity from public clouds. The approach of temporarily renting capacity to handle spikes in load is known as —cloud-bursting.

DESIRED FEATURES OF A CLOUD

Certain features of a cloud are essential to enable services that truly represent the cloud computing model and satisfy expectations of consumers, and cloud offerings must be having following features:

1. Self-service

2. Per-usage metered and billed
3. Elastic,
4. Customizable.

- **Self-Service**

Consumers of cloud computing services expect on-demand, nearly instant access to resources. To support this expectation, clouds must allow self-service access so that customers can request, customize, pay, and use services without intervention of human operators.

- **Per-Usage Metering and Billing**

Cloud computing eliminates up-front commitment by users, allowing them to request and use only the necessary amount. Services must be priced on a short-term basis (e.g., by the hour), allowing users to release (and not pay for) resources as soon as they are not needed. For these reasons, clouds must implement features to allow efficient trading of service such as pricing, accounting, and billing. Metering should be done accordingly for different types of service (e.g., storage, processing, and bandwidth) and usage promptly reported, thus providing greater transparency.

- **Elasticity**

Cloud computing gives the illusion of infinite computing resources available on demand. Therefore, users expect clouds to rapidly provide resources in any quantity at any time. In particular, it is expected that the additional resources can be (a) provisioned, possibly automatically, when an application load increases and (b) released when load decreases (scale up and down).

- **Customization**

In a multi-tenant cloud a great disparity between user needs is often the case. Thus, resources rented from the cloud must be highly customizable. In the case of infrastructure services, customization means allowing users to deploy specialized virtual appliances and to be given privileged (root) access to the virtual servers. Other service classes (PaaS and SaaS) offer less flexibility and are not suitable for general-purpose computing, but still are expected to provide a certain level of customization.

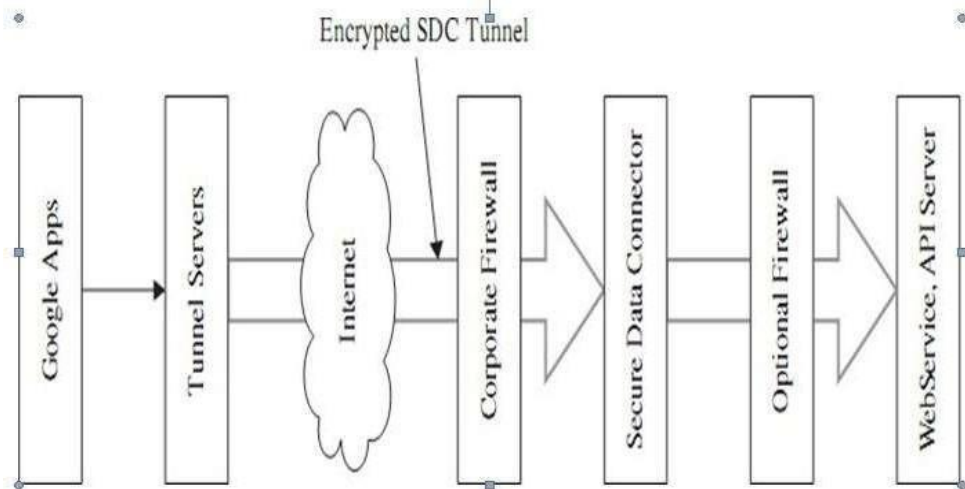
Google APP Engine

1. The app engine is a Cloud-based platform, is quite comprehensive and combines infrastructure as a service (IaaS), platform as a service (PaaS) and software as a service (SaaS). The app engine supports the delivery, testing and development of software on demand in a Cloud computing environment that supports millions of users and is highly scalable.
2. The company extends its platform and infrastructure to the Cloud through its app engine. It presents the platform to those who want to develop SaaS solutions at competitive costs.

3. It is a platform for hosting web applications in Google managed data centers. It is cloud-computing technology which virtualizes applications across multiple servers and data centers. Running your web application in Google infrastructure and Support different runtime environments

Java (JRE 6 with limitation, Servlet 2.5, JDO, JPA) Python (2.5.2)

1. Apps run in sandbox.
2. Automatic scaling and load balancing
3. No server restart, no network issues



1. The SDC constructs an encrypted connection between the data source and Google Apps. As long as the data source is in the Google Apps domain to the Google tunnel protocol servers, when the user wants to get the data, he/she will first send an authorized data requests to Google Apps, which forwards the request to the tunnel server.
2. The tunnel servers validate the request identity. If the identity is valid, the tunnel protocol allows the SDC to set up a connection, authenticate, and encrypt the data that flows across the Internet. At the same time, the SDC uses resource rules to validate whether a user is authorized to access a specified resource.
3. When the request is valid, the SDC performs a network request. The server validates the signed request, checks the credentials, and returns the data if the user is authorized. From the perspective of cloud storage services, data integrity depends on the security of operations while in storage in addition to the security of the uploading and downloading sessions. The uploading session can only ensure that the data received by the cloud storage is the data that the user uploaded; The downloading session can guarantee the data that the user retrieved is the data cloud storage recorded. Unfortunately, this procedure applied on cloud storage services cannot guarantee data integrity. First, assume that Alice, a company CFO, stores the company financial data at a cloud storage service provided by Eve. And then Bob, the company administration chairman, downloads the data from the cloud.

There are three important concerns in this simple procedure:

1. **Confidentiality.** Eve is considered as an untrustworthy third party, Alice and Bob do not want reveal the data to Eve.
2. **Integrity.** As the administrator of the storage service, Eve has the capability to play with the data in hand. How can Bob be confident that the data he fetched from Eve are the same as what was sent by Alice? Are there any measures to guarantee that the data have not been tampered by Eve?
3. **Repudiation.** If Bob finds that the data have been tampered with, is there any evidence for him to demonstrate that it is Eve who should be responsible for the fault? Similarly, Eve also needs certain evidence to prove her innocence.

GoogleAPP Engine Solutions for Missing Link

1. Third Authority Certified(TAC)
2. Secret Key Sharing(SKS)

Four Solutions

1. Neither TAC nor SKS
2. With SKS but without TAC
3. With TAC but without SKS O With Both TAC and SKS

Google is a leader in web-based applications.

Google is a leader in web-based applications and leading searching engine in the world. so it's not surprising that the company also offers cloud development services.

1. These services come in the form of the Google App Engine, which enables developers to build their own web applications utilizing the same infrastructure that powers Google's powerful applications.
2. The Google App Engine provides a fully integrated application environment. Using Google's development tools and computing cloud, App Engine applications are easy to build, easy to maintain, and easy to scale. All you have to do

Features of App Engine

1. These are covered by the depreciation policy and the service-level agreement of the app engine. Any changes made to such a feature are backward-compatible and implementation of such a feature is usually stable. These include data storage, retrieval, and search; communications; process management; computation; app configuration and management.
2. Data storage, retrieval, and search include features such as HRD migration tool, Google Cloud SQL, logs, datastore, dedicated Memcache, blobstore, Memcache and search.
3. Communications include features such as XMPP. channel, URL fetch, mail, and Google Cloud Endpoints.
4. Process management includes features like scheduled tasks and task queue. Computation includes images.

5. App management and configuration cover app identity, users, capabilities, traffic splitting, modules, SSL for custom domains, modules, remote access, and multi-tenancy.

Centralizing email Communications

1. The key here is to enable anywhere/anytime access to email. Precloud computing, your email access was via a single computer, which also stored all your email messages.
2. For this purpose, you probably used a program like Microsoft Outlook or Outlook Express, installed on your home computer.
3. To check your home email from work, it took a bit of juggling and perhaps the use of your ISP's email access web page. That web page was never in sync with the messages on your home PC, of course, which is just the start of the problems with trying to communicate in this fashion.
4. A better approach is to use a web-based email service, such as Google's Gmail (mail.google.com), Microsoft's Windows Live Hotmail (mail.live.com), or Yahoo! Mail (mail.yahoo.com).
5. These services place your email inbox in the cloud and you can access it from any computer connected to the Internet.

Collaborating via Web-Based Communication Tools-GMAIL

1. Gmail offers a few unique features that set it apart from the web-based email crowd.
2. First, Gmail doesn't use folders. With Gmail you can't organize your mail into folders, as you can with the other services.
3. Instead, Gmail pushes the search paradigm as the way to find the messages you want—not a surprise, given Google's search-centric business model.
4. Gmail does, however, let you —tag each message with one or more labels. This has the effect of creating virtual folders, as you can search and sort your messages by any of their labels.
5. In addition, Gmail groups together related email messages in what Google calls conversations. Yahoo! Mail (mail.yahoo.com)
6. There is another web mail service, provided by the popular Yahoo! search site.
7. The basic Yahoo! Mail is free and can be accessed from any PC, using any web browser.
8. Yahoo! also offers a paid service called Yahoo! Mail Plus that lets you send larger messages and offers offline access to your messages via POP email clients.

Web Mail Services

1. AOL Mail (mail.aol.com)
2. BigString (www.bigstring.com)
3. Excite Mail (mail.excite.com)

4. FlashMail (www.flashmail.com)
5. GMX Mail (www.gmx.com)
6. Inbox.com (www.inbox.com)
7. Lycos Mail (mail.lycos.com)
8. Mail.com (www.mail.com)
9. Zoho Mail (zoho.mail.com)

An Introduction to Data Security:

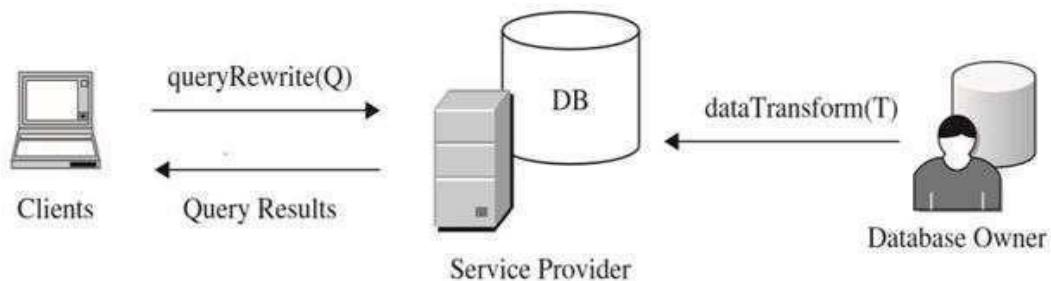
Data Security

1. Data Security defines as Information in a cloud environment has much more dynamism and fluidity than information that is static on a desktop or in a network folder
2. Nature of cloud computing dictates that data are fluid objects, accessible from a multitude of nodes and geographic locations and, as such, must have a data security methodology that takes this into account while ensuring that this fluidity is not compromised.
3. The idea of content-centric or information-centric protection, being an inherent part of a data object is a development out of the idea of the —de-perimeterization of the enterprise.
4. This idea was put forward by a group of Chief Information Officers (CIOs) who formed an organization called the Jericho Forum

TECHNOLOGIES FOR DATA SECURITY IN CLOUD COMPUTING

Unique issues of the cloud data storage platform from a few different perspectives

1. Database Outsourcing and Query Integrity Assurance
 - I. Storing data into and fetching data from devices and machines behind a cloud are essentially a novel form of database outsourcing



2. Data Integrity in Untrustworthy Storage
 - I. The fear of losing data or data corruption
 - II. Relieve the users' fear by providing technologies that enable users to check the

integrity of their data

3. Web-Application-Based Security

- I. Once the dataset is stored remotely, a Web browser is one of the most convenient approaches that end users can use to access their data on remote services
- II. Web security plays a more important role for cloud computing

1. Multimedia Data Security

1. With the development of high-speed network technologies and large bandwidth connections, more and more multimedia data are being stored and shared in cyber space
2. The security requirements for video, audio, pictures, or images are different from other applications

CLOUD COMPUTING AND IDENTITY

Digital identity

1. Digital identity holds the key to flexible data security within a cloud Environment
2. A digital identity represents who we are and how we interact with others on-line.
3. **Access, identity, and risk** are three variables that can become inherently connected when applied to the security of data, because access and risk are directly proportional: As access increases, so then risk to the security of the data increases.
4. Access controlled by identifying the actor attempting the access is the most logical manner of performing this operation.
5. Ultimately, digital identity holds the key to securing data, if that digital identity can be programmatically linked to security policies controlling the post-access usage of data.

Identity, Reputation, and Trust

1. Reputation is a real-world commodity; that is a basic requirement of human-to-human relationships
2. Our basic societal communication structure is built upon the idea of reputation and trust.
3. Reputation and its counter value, trust, is easily transferable to a digital realm:
4. EBay, for example, having partly built a successful business model on the strength of a ratings system, builds up the reputation of its buyers and sellers through successful (or unsuccessful) transactions.
5. These types of reputation systems can be extremely useful when used with a digital identity.
6. They can be used to associate varying levels of trust with that identity, which in turn can be used to define the level (granular variations) of security policy applied to data resources that the individual wishes to access

User-Centric Identity:

1. Digital identities are a mechanism for identifying an individual, particularly within a cloud environment and identity ownership being placed upon the individual is known as user- centric identity
2. It allows users to consent and control how their identity (and the individual identifiers making up the identity, the claims) is used.
3. This reversal of ownership away from centrally managed identity platforms (enterprise- centric) has many advantages.
4. This includes the potential to improve the privacy aspects of a digital identity, by giving an individual the ability to apply permission policies based on their identity and to control which aspects of that identity are divulged
5. An identity may be controllable by the end user, to the extent that the user can then decide what information is given to the party relying on the identity

Information Card:

- Information cards permit a user to present to a Web site or other service (relying party) one or more claims, in the form of a software token, which may be used to uniquely identify that user.
- They can be used in place of user name/ passwords, digital certificates, and other identification systems, when user identity needs to be established to control access to a Web site or other resource, or to permit digital signing
- Information cards are part of an identity meta-system consisting of:
 - **Identity providers (IdP)**, who provision and manage information cards with specific claims, to users.
 - **Users** who own and utilize the cards to gain access to Web sites and other resources that support information cards.
 - **An identity selector/service**, which is a piece of software on the user's desktop or in the cloud that allows a user to select and manage their cards.
 - **Relying parties.** These are the applications, services & so on, that can use an information card to authenticate a person and to then authorize an action such as logging onto a Web site, accessing a document, signing content, and so on.
 - Each information card is associated with a set of claims which can be used to identify the user. These claims include identifiers such as name, email address post code.

Using Information Cards to Protect Data

1. Information cards are built around a set of open standards devised by a consortium that includes Microsoft, IBM, Novell, and so on.
2. The original remit of the cards was to create a type of single sign on system for the Internet, to help users to move away from the need to remember multiple passwords.

3. However, the information card system can be used in many more ways.
4. Because an information card is a type of digital identity, it can be used in the same way that other digital identities can be used.
5. For example, an information card can be used to digitally sign data and content and to control access to data and content. One of the more sophisticated uses of an information card is the advantage given to the cards by way of the claims system.

Cloud Computing and Data Security Risk

1. Cloud computing is a development that is meant to allow more open accessibility and easier and improved data sharing.
2. Data are uploaded into a cloud and stored in a data center, for access by users from that data center; or in a more fully cloud-based model, the data themselves are created in the cloud and stored and accessed from the cloud (again via a data center).
3. The most obvious risk in this scenario is that associated with the storage of that data. A user uploading or creating cloud-based data include those data that are stored and maintained by a third-party cloud provider such as Google, Amazon, Microsoft, and so on.

This action has several risks associated with it:

- i. Firstly, it is necessary to protect the data during upload into the data center to ensure that the data do not get hijacked on the way into the database.
- ii. Secondly, it is necessary to store the data in the data center to ensure that they are encrypted at all times.
- iii. Thirdly, and perhaps less obvious, the access to those data need to be controlled; this control should also be applied to the hosting company, including the administrators of the data center.
- iv. In addition, an area often forgotten in the application of security to a data resource is the protection of that resource during its use

Data security risks are compounded by the open nature of cloud computing.

1. Access control becomes a much more fundamental issue in cloud-based systems because of the accessibility of the data
2. Information-centric access control (as opposed to access control lists) can help to balance improved accessibility with risk, by associating access rules with different data objects within an open and accessible platform, without losing the inherent usability of that platform
3. A further area of risk associated not only with cloud computing, but also with traditional network computing, is the use of content after access.

4. The risk is potentially higher in a cloud network, for the simple reason that the information is outside of your corporate walls

Data-centric mashups

1. that are used to perform business processes around data creation and dissemination—by their very nature, can be used to hijack data, leaking sensitive information and/or affecting integrity of that data
2. Cloud computing, more than any other form of digital communication technology, has created a need to ensure that protection is applied at the inception of the information, in a content centric manner, ensuring that a security policy becomes an integral part of that data throughout its life cycle.

Encryption

1. It is a vital component of the protection policy, but further controls over the access of that data and on the use of the data must be met.
2. In the case of mashups, the controlling of access to data resources, can help to alleviate the security concerns by ensuring that mashup access is authenticated.
3. Linking security policies, as applied to the use of content, to the access control method offer a way of continuing protection of data, post access and throughout the life cycle; this type of data security philosophy must be incorporated into the use of cloud computing to alleviate security risks.

UNIT- V

SLA Management in cloud computing: Traditional Approaches to SLO Management, Types of SLA, Life Cycle of SLA, SLA Management in Cloud.

SLA MANAGEMENT IN CLOUD COMPUTING

In the early days of web-application deployment, performance of the application at peak load was a single important criterion for provisioning server resources. Provisioning in those days involved deciding hardware configuration, determining the number of physical machines, and acquiring them upfront so that the overall business objectives could be achieved. The web applications were hosted on these dedicated individual servers within enterprises' own server rooms. These web applications were used to provide different kinds of e-services to various clients. Typically, the service-level objectives (SLOs) for these applications were response time and throughput of the application end-user requests. The capacity buildup was to cater to the estimated peak load experienced by the application. The activity of determining the number of servers and their capacity that could satisfactorily serve the application end-user requests at peak loads is called capacity planning. An example scenario where two web applications, application A and application B, are hosted on a separate set of dedicated servers within the enterprise-owned server rooms is shown in Figure 16.1. The planned capacity for each of the applications to run successfully is three servers. As the number of web applications grew, the server rooms in the organization became large and such server rooms were known as data centers. These data centers were owned and managed by the enterprises themselves.

414 SLA MANAGEMENT IN CLOUD COMPUTING: A SERVICE PROVIDER'S PERSPECTIVE

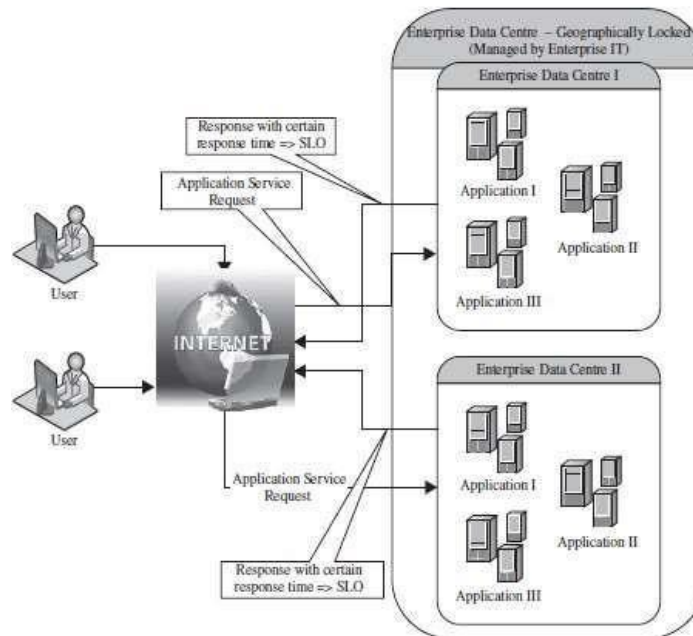


FIGURE 16.1. Hosting of applications on servers within enterprise's data centers.

TRADITIONAL APPROACHES TO SLO MANAGEMENT

Traditionally, load balancing techniques and admission control mechanisms have been used to provide guaranteed quality of service (QoS) for hosted web applications. These mechanisms can be viewed as the first attempt towards managing the SLOs. In the following subsections we discuss the existing approaches for load balancing and admission control for ensuring QoS.

Load Balancing The objective of a load balancing is to distribute the incoming requests onto a set of physical machines, each hosting a replica of an application, so that the load on the machines is equally distributed. The load balancing algorithm executes on a physical machine that interfaces with the clients. This physical machine, also called the front-end node, receives the incoming requests and distributes these requests to different physical machines for further execution.

This set of physical machines is responsible for serving the incoming requests and are known as the back-end nodes.

TYPES OF SLA:

Service-level agreement provides a framework within which both seller and buyer of a service can pursue a profitable service business relationship. It outlines the broad understanding between the service provider and the service consumer for conducting business and forms the basis for maintaining a mutually beneficial relationship. From a legal perspective, the necessary terms and conditions that bind the service provider to provide services continually to the service consumer are formally defined in SLA.

SLA can be modeled using web service-level agreement (WSLA) language specification. Although WSLA is intended for web-service-based applications, it is equally applicable for hosting of applications. Service-level parameter, metric, function, measurement directive, service-level objective, and penalty are some of the important

TABLE 16.1. Key Components of a Service-Level Agreement

Service-Level Parameter	Describes an observable property of a service whose value is measurable.
Metrics	These are definitions of values of service properties that are measured from a service-providing system or computed from other metrics and constants. Metrics are the key instrument to describe exactly what SLA parameters mean by specifying how to measure or compute the parameter values.
Function	A function specifies how to compute a metric's value from the values of other metrics and constants. Functions are central to describing exactly how SLA parameters are computed from resource metrics.
Measurement directives	These specify how to measure a metric.

components of WSLA and are described in Table 16.1.

There are two types of SLAs from the perspective of application hosting. These are described in detail here.

Infrastructure SLA: The infrastructure provider manages and offers guarantees on availability of the infrastructure, namely, server machine, power, network connectivity, and so on. Enterprises manage themselves, their applications that are deployed on these server machines. The machines are leased to the customers and are isolated from machines of other customers. In such dedicated hosting environments, a practical example of service-level guarantees offered by infrastructure providers is shown in Table 16.2.

Application SLA: In the application co-location hosting model, the server capacity is available to the applications based solely on their resource demands. Hence, the service providers are flexible in allocating and de-allocating computing resources among the co-located applications.

Therefore, the service providers are also responsible for ensuring to meet their customer's application SLOs. For example, an enterprise can have the following application SLA with a service provider for one of its application.

LIFE CYCLE OF SLA:

Each SLA goes through a sequence of steps starting from identification of terms and conditions, activation and monitoring of the stated terms and conditions, and eventual termination of contract once the hosting relationship ceases to exist. Such a sequence of steps is called SLA life cycle and consists of the following five phases:

1. Contract definition
2. Publishing and discovery
3. Negotiation
4. Operationalization
5. De-commissioning

Here, we explain in detail each of these phases of SLA life cycle.

Contract Definition: Generally, service providers define a set of service offerings and corresponding SLAs using standard templates. These service offerings form a catalog. Individual SLAs for enterprises can be derived by customizing these base SLA templates.

Publication and Discovery: Service provider advertises these base service offerings through standard publication media, and the customers should be able to locate the service provider by searching the catalog. The customers can search different competitive offerings and shortlist a few that fulfill their requirements for further negotiation.

Negotiation: Once the customer has discovered a service provider who can meet their application hosting need, the SLA terms and conditions needs to be mutually agreed upon before signing the agreement for hosting the application. For a standard packaged application

which is offered as service, this phase could be automated. For customized applications that are hosted on cloud platforms, this phase is manual. The service provider needs to analyze the application's behavior with respect to scalability and performance before agreeing on the specification of SLA. At the end of this phase, the SLA is mutually agreed by both customer and provider and is eventually signed off. SLA negotiation can utilize the WS-negotiation specification.

Operationalization: SLA operation consists of SLA monitoring, SLA accounting, and SLA enforcement. SLA monitoring involves measuring parameter values and calculating the metrics defined as a part of SLA and determining the deviations. On identifying the deviations, the concerned parties are notified. SLA accounting involves capturing and archiving the SLA adherence for compliance.

As part of accounting, the application's actual performance and the performance guaranteed as a part of SLA is reported. Apart from the frequency and the duration of the SLA breach, it should also provide the penalties paid for each SLA violation. SLA enforcement involves taking appropriate action when the runtime monitoring detects a SLA violation. Such actions could be notifying the concerned parties, charging the penalties besides other things. The different policies can be expressed using a subset of the Common Information Model (CIM) [9]. The CIM model is an open standard that allows expressing managed elements of data center via relationships and common objects.

De-commissioning: SLA decommissioning involves termination of all activities performed under a particular SLA when the hosting relationship between the service provider and the service consumer has ended. SLA specifies the terms and conditions of contract termination and specifies situations under which the relationship between a service provider and a service consumer can be considered to be legally ended.

SLA MANAGEMENT IN CLOUD: SLA management of applications hosted on cloud platforms involves five phases.

1. Feasibility
2. On-boarding
3. Pre-production
4. Production
5. Termination

Different activities performed under each of these phases are shown in Figure 16.7. These activities are explained in detail in the following subsections.

Feasibility Analysis:

MSP conducts the feasibility study of hosting an application on their cloud platforms. This study involves three kinds of feasibility: (1) technical feasibility, (2) infrastructure feasibility, and (3) financial feasibility. The technical feasibility of an application implies determining the following:

1. Ability of an application to scale out.
2. Compatibility of the application with the cloud platform being used within the MSP's data center.
3. The need and availability of a specific hardware and software required for hosting and running of the application.
4. Preliminary information about the application performance and whether they can be met by the MSP.



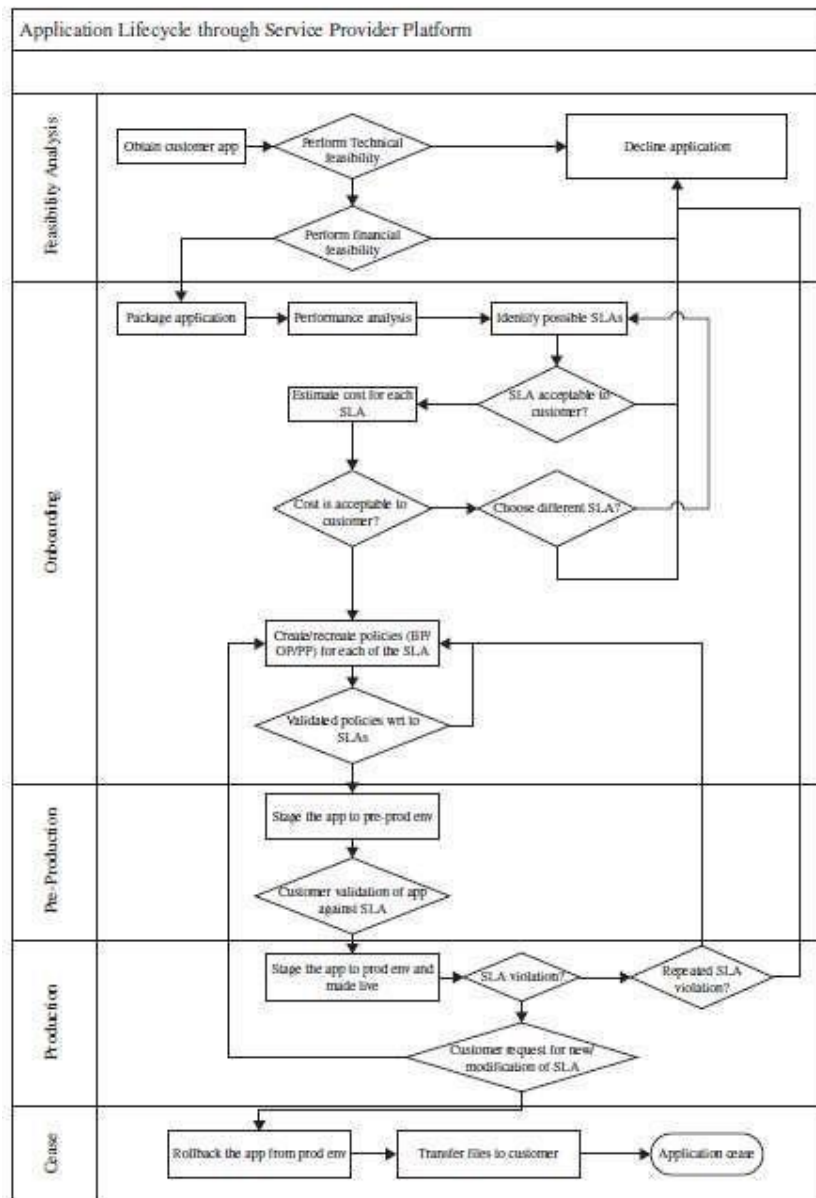


FIGURE 16.7. Flowchart of the SLA management in cloud.

On-Boarding of Application

Once the customer and the MSP agree in principle to host the application based on the findings of the feasibility study, the application is moved from the customer servers to the hosting platform. Moving an application to the MSP's hosting platform is called on-boarding. As part of the on-boarding activity, the MSP understands the application runtime characteristics using runtime profilers. This helps the MSP to identify the possible SLAs that can be offered to the customer for that application. This also helps in creation of the necessary policies (also called rule sets) required to guarantee the SLOs mentioned in the

application SLA. The application is accessible to its end users only after the on-boarding activity is completed.

Preproduction

Once the determination of policies is completed as discussed in previous phase, the application is hosted in a simulated production environment. It facilitates the customer to verify and validate the MSP's findings on application's runtime characteristics and agree on the defined SLA. Once both parties agree on the cost and the terms and conditions of the SLA, the customer sign-off is obtained. On successful completion of this phase the MSP allows the application to go on-live.

Production

In this phase, the application is made accessible to its end users under the agreed SLA. However, there could be situations when the managed application tends to behave differently in a production environment compared to the preproduction environment. This in turn may cause sustained breach of the terms and conditions mentioned in the SLA. Additionally, customer may request the MSP for inclusion of new terms and conditions in the SLA. If the application SLA is breached frequently or if the customer requests for a new non-agreed SLA, the on-boarding process is performed again. In the case of the former, on-boarding activity is repeated to analyze the application and its policies with respect to SLA fulfillment. In case of the latter, a new set of policies are formulated to meet the fresh terms and conditions of the SLA.

Termination

When the customer wishes to withdraw the hosted application and does not wish to continue to avail the services of the MSP for managing the hosting of its application, the termination activity is initiated. On initiation of termination, all data related to the application are transferred to the customer and only the essential information is retained for legal compliance. This ends the hosting relationship between the two parties for that application, and the customer sign-off is obtained.

Terminology

Basic Terms

- **Cloud Computing:** Delivery of computing services (servers, storage, databases, networking, software) over the internet (“the cloud”).
- **Public Cloud:** Cloud services offered over the internet and shared across multiple organizations (e.g., AWS, Azure, Google Cloud).
- **Private Cloud:** Cloud infrastructure dedicated to a single organization, either hosted on-premises or by a provider.
- **Hybrid Cloud:** A mix of public and private clouds, enabling data and applications to be shared between them.
- **Multi-Cloud:** Using multiple cloud providers for different services or applications.

Cloud Service Models

- **IaaS (Infrastructure as a Service):** Virtualized computing resources over the internet (e.g., Amazon EC2, Azure VM).
- **PaaS (Platform as a Service):** Platform allowing customers to develop, run, and manage applications (e.g., Google App Engine).
- **SaaS (Software as a Service):** Software delivered via the cloud and accessed through a browser (e.g., Dropbox, Salesforce).

Deployment & Management

- **Virtual Machine (VM):** Emulation of a physical computer running an OS, used in IaaS.
- **Container:** Lightweight, standalone software package that includes code and dependencies (e.g., Docker).
- **Orchestration:** Automated arrangement, coordination, and management of containers (e.g., Kubernetes).
- **Serverless Computing:** Running code without managing servers; automatically scales (e.g., AWS Lambda).
- **Auto-scaling:** Automatic adjustment of cloud resources based on load.

Storage & Networking

- **Object Storage:** Storage for unstructured data, often used for media, backups (e.g., Amazon S3).
- **Block Storage:** Data stored in blocks, used for databases and VM file systems.
- **Virtual Private Cloud (VPC):** Isolated network within a public cloud.
- **Content Delivery Network (CDN):** Network of servers that deliver content based on user location (e.g., Cloudflare, AWS CloudFront).

Security & Compliance

- **IAM (Identity and Access Management):** Framework for managing digital identities and access.
- **Encryption:** Converting data into a secure format.
- **Compliance:** Adhering to regulations (e.g., GDPR, HIPAA).
- **Firewall:** Network security system that monitors and controls incoming and outgoing traffic.

Billing & Monitoring

- **Pay-as-you-go:** Billing based on actual usage.
- **Reserved Instances:** Prepaid cloud instances for lower cost.
- **Cloud Monitoring:** Tracking performance, usage, and uptime (e.g., CloudWatch, Azure Monitor).