

1.1 Introduction

Cloud computing is a transformative computing paradigm that involves delivering applications and services over the internet. Many of the underlying technologies that are the foundation of cloud computing have existed for quite some time. Cloud computing involves provisioning of computing, networking and storage resources on demand and providing these resources as metered services to the users, in a "pay as you go" model. In this chapter you will learn about the various deployment models, service models, characteristics, driving factors and challenges of cloud computing.

1.1.1 Definition of Cloud Computing

The U.S. National Institute of Standards and Technology (NIST) defines cloud computing as [1]:

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

1.2 Characteristics of Cloud Computing

NIST further identifies five essential characteristics of cloud computing:

On-demand self service

Cloud computing resources can be provisioned on-demand by the users, without requiring interactions with the cloud service provider. The process of provisioning resources is automated.

Broad network access

Cloud computing resources can be accessed over the network using standard access mechanisms that provide platform-independent access through the use of heterogeneous client platforms such as workstations, laptops, tablets and smartphones.

Resource pooling

The computing and storage resources provided by cloud service providers are pooled to serve multiple users using multi-tenancy. Multi-tenant aspects of the cloud allow multiple users to be served by the same physical hardware. Users are assigned virtual resources that run on top of the physical resources. Various forms of virtualization approaches such as full virtualization, para-virtualization and hardware virtualization are described in Chapter 2.

Rapid elasticity

Cloud computing resources can be provisioned rapidly and elastically. Cloud resources can be rapidly scaled up or down based on demand. Two types of scaling options exist:

- **Horizontal Scaling (scaling out):** Horizontal scaling or scaling-out involves launching and provisioning additional server resources.

- **Vertical Scaling (scaling up):** Vertical scaling or scaling-up involves changing the computing capacity assigned to the server resources while keeping the number of server resources constant.

Measured service

Cloud computing resources are provided to users on a pay-per-use model. The usage of the cloud resources is measured and the user is charged based on some specific metric. Metrics such as amount of CPU cycles used, amount of storage space used, number of network I/O requests, etc. are used to calculate the usage charges for the cloud resources.

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

Performance

Cloud computing provides improved performance for applications since the resources available to the applications can be scaled up or down based on the dynamic application workloads.

Reduced costs

Cloud computing provides cost benefits for applications as only as much computing and storage resources as required can be provisioned dynamically, and upfront investment in purchase of computing assets to cover worst case requirements is avoided. This saves significant cost for organizations and individuals. Applications can experience large variations in the workloads which can be due to seasonal or other factors. For example, e-Commerce applications typically experience higher workloads in holiday seasons. To ensure market readiness of such applications, adequate resources need to be provisioned so that the applications can meet the demands of specified workload levels and at the same time ensure that service level agreements are met.

Outsourced Management

Cloud computing allows the users (individuals, large organizations, small and medium enterprises and governments) to outsource the IT infrastructure requirements to external cloud providers. Thus, the consumers can save large upfront capital expenditures in setting up the IT infrastructure and pay only for the operational expenses for the cloud resources used. The outsourced nature of the cloud services provides a reduction in the IT infrastructure management costs.

Reliability

Applications deployed in cloud computing environments generally have a higher reliability since the underlying IT infrastructure is professionally managed by the cloud service. Cloud service providers specify and guarantee the reliability and availability levels for their cloud resources in the form of service level agreements (SLAs). Most cloud providers promise 99.99% uptime guarantee for the cloud resources, which may often be expensive to achieve with in-house IT infrastructure.

Multi-tenancy

The multi-tenanted approach of the cloud allows multiple users to make use of the same shared resources. Modern applications such as e-Commerce, Business-to-Business, Banking and

Financial, Retail and Social Networking applications that are deployed in cloud computing environments are multi-tenanted applications. Multi-tenancy can be of different forms:

- **Virtual multi-tenancy:** In virtual multi-tenancy, computing and storage resources are shared among multiple users. Multiple tenants are served from virtual machines (VMs) that execute concurrently on top of the same computing and storage resources.
- **Organic multi-tenancy:** In organic multi-tenancy every component in the system architecture is shared among multiple tenants, including hardware, OS, database servers, application servers, load balancers, etc. Organic multi-tenancy exists when explicit multi-tenant design patterns are coded into the application.

1.3 Cloud Models

1.3.1 Service Models

Cloud computing services are offered to users in different forms. NIST defines at least three cloud service models as follows:

Infrastructure-as-a-Service (IaaS)

IaaS provides the users the capability to provision computing and storage resources. These resources are provided to the users as virtual machine instances and virtual storage. Users can start, stop, configure and manage the virtual machine instances and virtual storage. Users can deploy operating systems and applications of their choice on the virtual resources provisioned in the cloud. The cloud service provider manages the underlying infrastructure. Virtual resources provisioned by the users are billed based on a pay-per-use paradigm. Common metering metrics used are the number of virtual machine hours used and/or the amount of storage space provisioned.

Platform-as-a-Service (PaaS)

PaaS provides the users the capability to develop and deploy application in the cloud using the development tools, application programming interfaces (APIs), software libraries and services provided by the cloud service provider. The cloud service provider manages the underlying cloud infrastructure including servers, network, operating systems and storage. The users, themselves, are responsible for developing, deploying, configuring and managing applications on the cloud infrastructure.

Software-as-a-Service (SaaS)

SaaS provides the users a complete software application or the user interface to the application itself. The cloud service provider manages the underlying cloud infrastructure including servers, network, operating systems, storage and application software, and the user is unaware of the underlying architecture of the cloud. Applications are provided to the user through a thin client interface (e.g., a browser). SaaS applications are platform independent and can be accessed from various client devices such as workstations, laptop, tablets and smartphones, running different operating systems. Since the cloud service provider manages both the application and data, the users are able to access the applications from anywhere.

Figure 1.1 shows the cloud computing service models and Figure 1.2 lists the benefits, characteristics and adoption of IaaS, PaaS and SaaS.

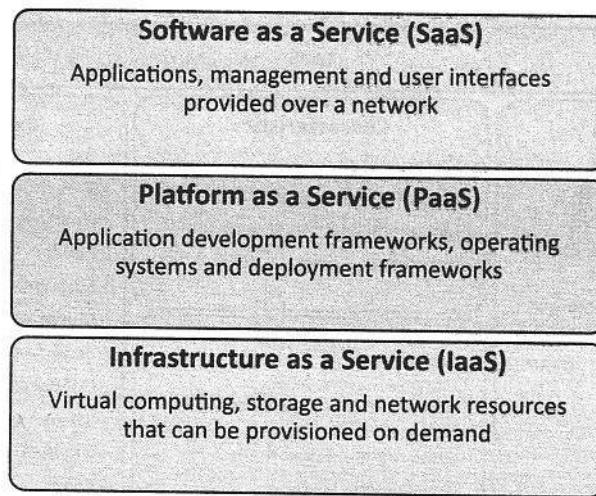


Figure 1.1: Cloud computing service models

1.3.2 Deployment Models

NIST also defines four cloud deployment models as follows:

Public cloud

In the public cloud deployment model, cloud services are available to the general public or a large group of companies. The cloud resources are shared among different users (individuals, large organizations, small and medium enterprises and governments). The cloud services are provided by a third-party cloud provider. Public clouds are best suited for users who want to use cloud infrastructure for development and testing of applications and host applications in the cloud to serve large workloads, without upfront investments in IT infrastructure.

Private cloud

In the private cloud deployment model, cloud infrastructure is operated for exclusive use of a single organization. Private cloud services are dedicated for a single organization. Cloud infrastructure can be setup on premise or off-premise and may be managed internally or by a third-party. Private clouds are best suited for applications where security is very important and organizations that want to have very tight control over their data.

Hybrid cloud

The hybrid cloud deployment model combines the services of multiple clouds (private or public). The individual clouds retain their unique identities but are bound by standardized or proprietary technology that enables data and application portability. Hybrid clouds are best suited for organizations that want to take advantage of secured application and data hosting on a private cloud, and at the same time benefit from cost savings by hosting shared applications and data in public clouds.

Community cloud

In the community cloud deployment model, the cloud services are shared by several organizations that have the same policy and compliance considerations. Community clouds are

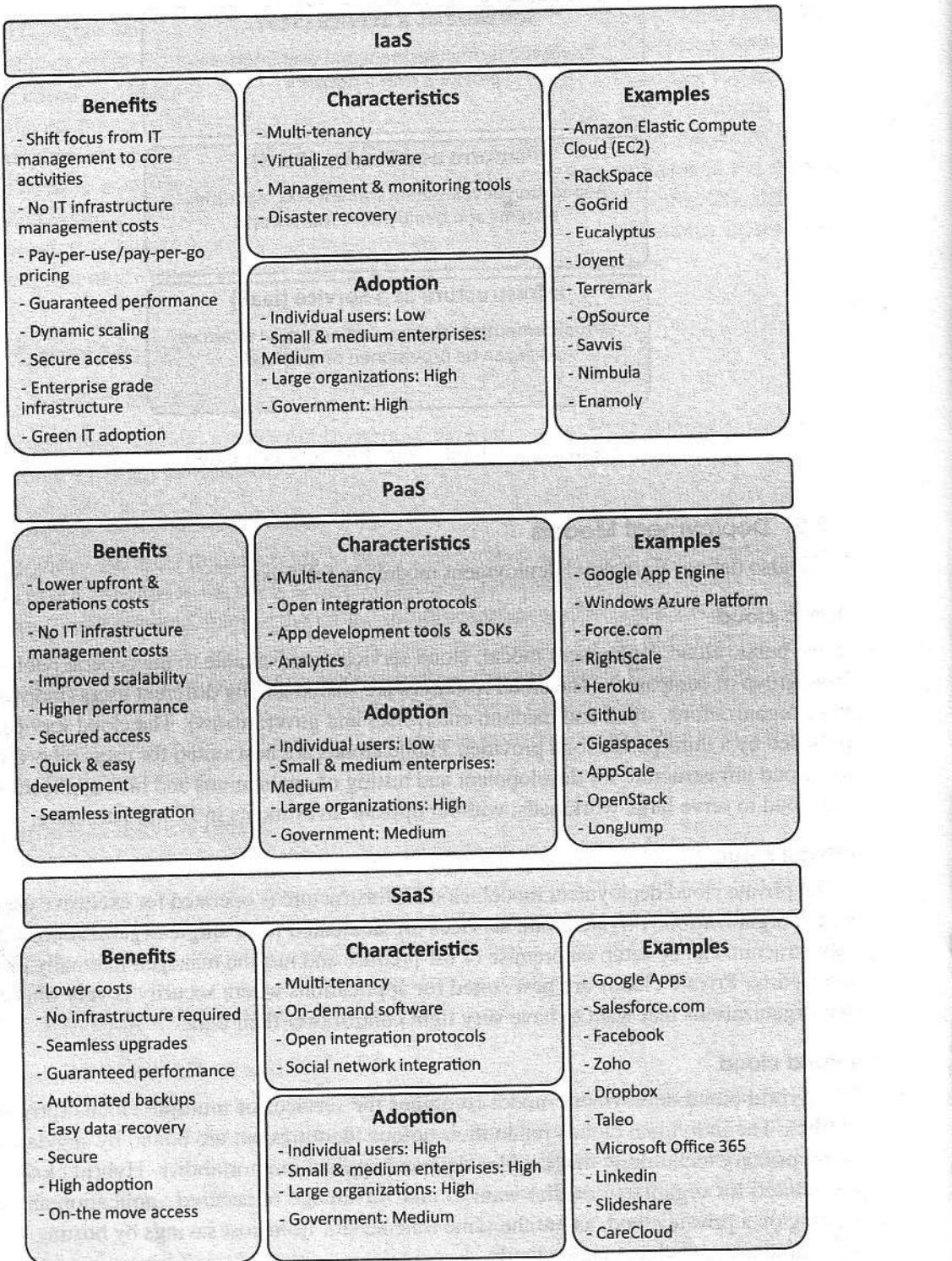


Figure 1.2: Benefits, characteristics and adoption of IaaS, PaaS and SaaS

best suited for organizations that want access to the same applications and data, and want the cloud costs to be shared with the larger group.

Figures 1.3 and 1.4 show the cloud deployment models.

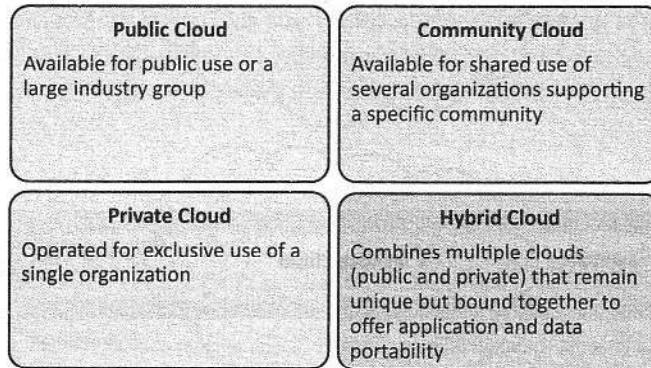


Figure 1.3: Cloud deployment models

1.4 Cloud Services Examples

1.4.1 IaaS: Amazon EC2, Google Compute Engine, Azure VMs

Amazon Elastic Compute Cloud (EC2) [3] is an Infrastructure as a Service (IaaS) offering from Amazon.com. EC2 (TM) is a web service that provides computing capacity in the form of virtual machines that are launched in Amazon's cloud computing environment. Amazon EC2 allows users to launch instances on demand using a simple web-based interface. Amazon provides pre-configured Amazon Machine Images (AMIs) which are templates of cloud instances. Users can also create their own AMIs with custom applications, libraries and data. Instances can be launched with a variety of operating systems. Users can load their applications on running instances and rapidly and easily increase or decrease capacity to meet the dynamic application performance requirements. With EC2, users can even provision hundreds or thousands of server instances simultaneously, manage network access permissions, and monitor usage resources through a web interface. Amazon EC2 provides instances of various computing capacities ranging from small instances (e.g., 1 virtual core with 1 EC2 compute unit, 1.7GB memory and 160GB instance storage) to extra large instances (e.g., 4 virtual cores with 2 EC2 compute units each, 15GB memory and 1690 GB instance storage). Amazon EC2 also provides instances with high memory, high CPU resources, cluster compute instances, cluster graphical processor unit (GPU) instances and high Input/Output (I/O) instances. The pricing model for EC2 instances is based on a pay-per-use model. Users are billed based on the number of instance hours used for on-demand instances. EC2 provides the option of reserving instances by one-time payment for each instance that the user wants to reserve. In addition to these on-demand and reserved instances, EC2 also provides spot instances that allow users to bid on unused Amazon EC2 capacity and run those instances for as long as their bid exceeds the current spot price. Amazon EC2 provides a number of powerful

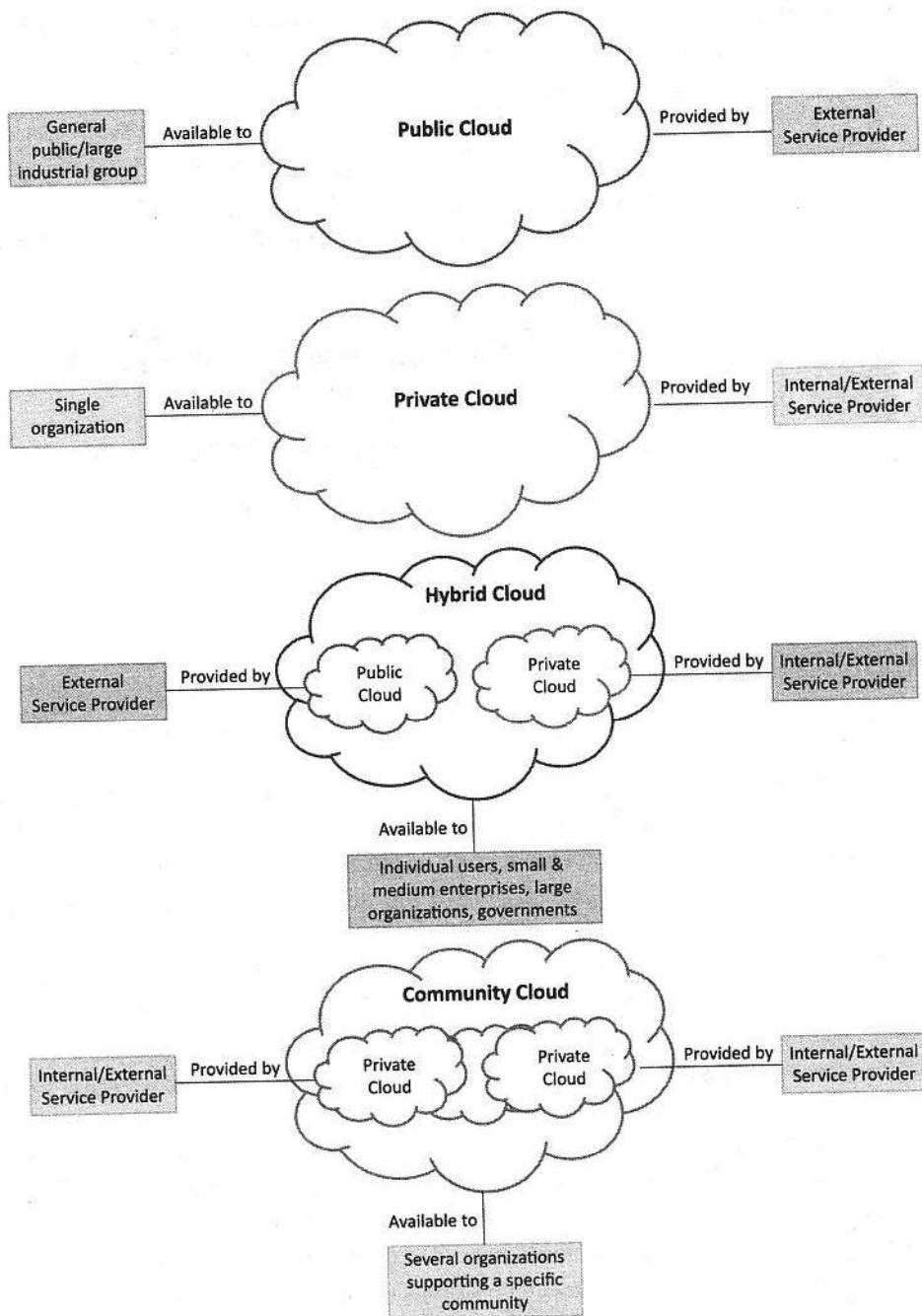


Figure 1.4: Cloud deployment models

features for building scalable and reliable applications such as auto scaling and elastic load balancing. Figure 1.5 shows a screenshot of Amazon EC2 dashboard.

Google Compute Engine (GCE) [4] is an IaaS offering from Google. GCE provides virtual machines of various computing capacities ranging from small instances (e.g., 1

virtual core with 1.38 GCE unit and 1.7GB memory) to high memory machine types (e.g., 8 virtual cores with 22 GCE units and 52GB memory). Figure 1.6 shows a screenshot of Google Compute Engine dashboard.

Windows Azure Virtual Machines [83] is an IaaS offering from Microsoft. Azure VMs provides virtual machines of various computing capacities ranging from small instances (1 virtual core with 1.75GB memory) to memory intensive machine types (8 virtual cores with 56GB memory). Figure 1.7 shows a screenshot of Windows Azure Virtual Machines dashboard.

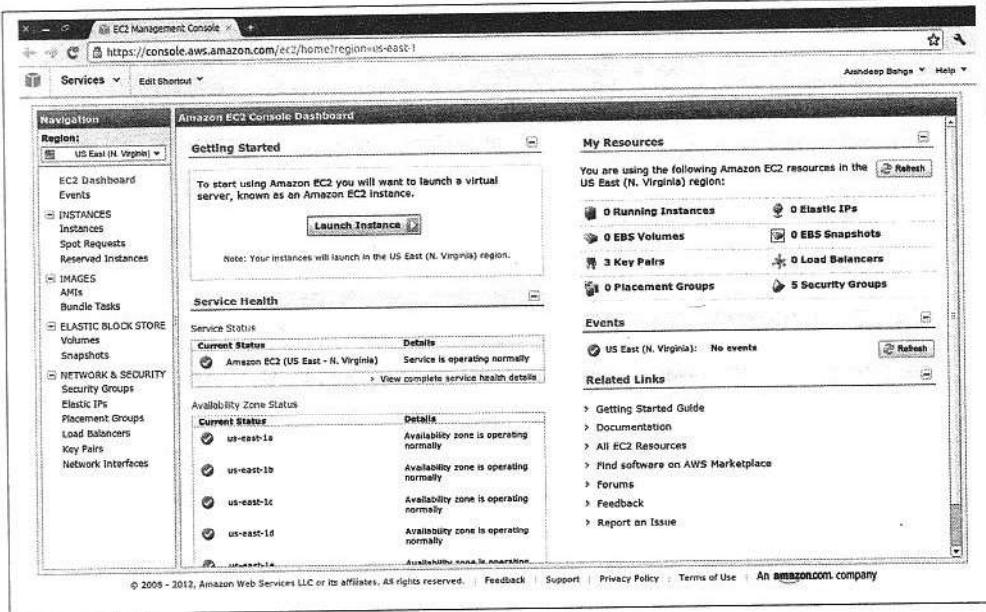


Figure 1.5: Amazon EC2 dashboard

1.4.2 PaaS: Google App Engine

Google App Engine (GAE) [105] is a Platform-as-a-Service (PaaS) offering from Google. GAE(TM) is a cloud-based web service for hosting web applications and storing data. GAE allows users to build scalable and reliable applications that run on the same systems that power Google's own applications. GAE provides a software development kit (SDK) for developing web applications software that can be deployed on GAE. Developers can develop and test their applications with GAE SDK on a local machine and then upload it to GAE with a simple click of a button. Applications hosted in GAE are easy to build, maintain and scale. Users don't need to worry about launching additional computing instances when the application load increases. GAE provides seamless scalability by launching additional instances when application load increases. GAE provides dynamic web serving based on common web technologies. Applications hosted in GAE can use dynamic technologies. GAE provides automatic scaling and load balancing capability.

GAE supports applications written in several programming languages. With GAE's

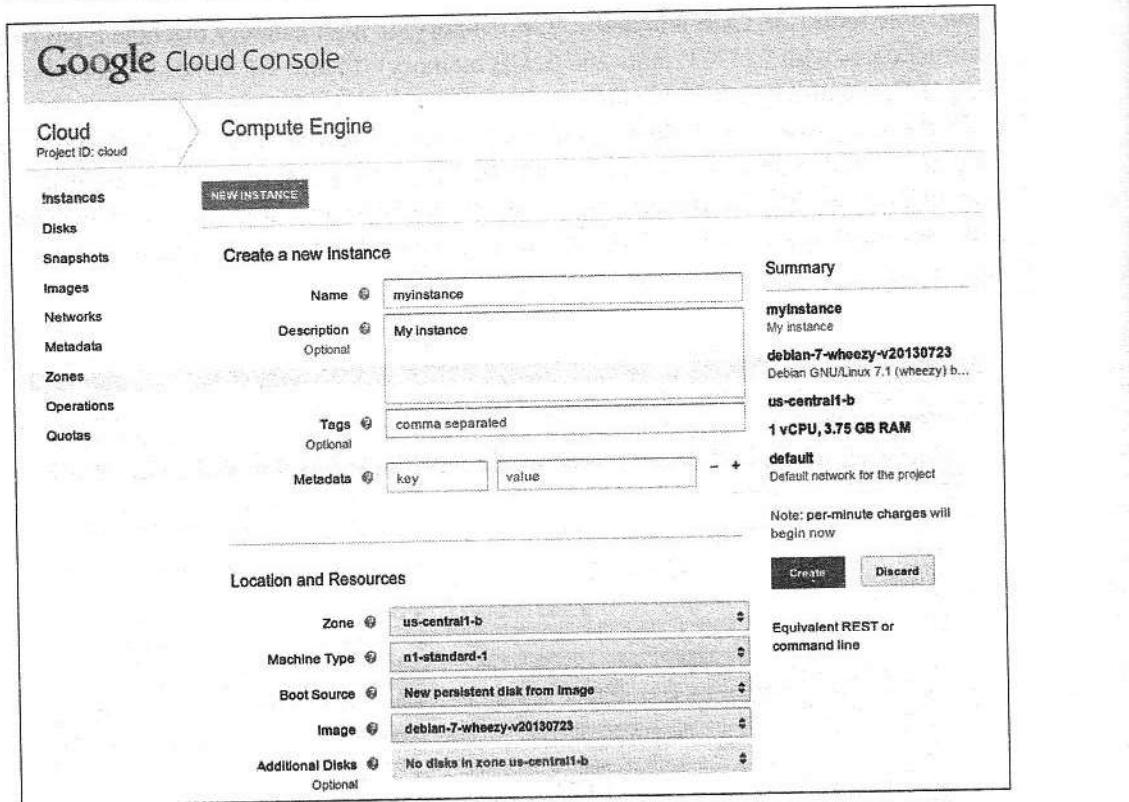


Figure 1.6: Google Compute Engine dashboard

Java runtime environment developers can build applications using Java programming language and standard Java technologies such as Java Servlets. GAE also provides runtime environments for Python and Go programming languages. Applications hosted in GAE run in secure sandbox with limited access to the underlying operating system and hardware. The benefit of hosting applications in separate sandboxes is that GAE can distribute web requests for applications across multiple servers thus providing scalability and security.

The pricing model for GAE is based on the amount of computing resources used. GAE provides free computing resources for applications up to a certain limit. Beyond that limit, users are billed based on the amount of computing resources used, such as amount bandwidth consumed, number of resources instance hours for front-end and back-end instances, amount of stored data, channels, and recipients emailed. Figure 1.8 shows a screenshot of GAE dashboard.

1.4.3 SaaS: Salesforce

Salesforce [7] Sales Cloud(TM) is a cloud-based customer relationship management (CRM) Software-as-a-Service (SaaS) offering. Users can access CRM application from anywhere through internet-enabled devices such as workstations, laptops, tablets and

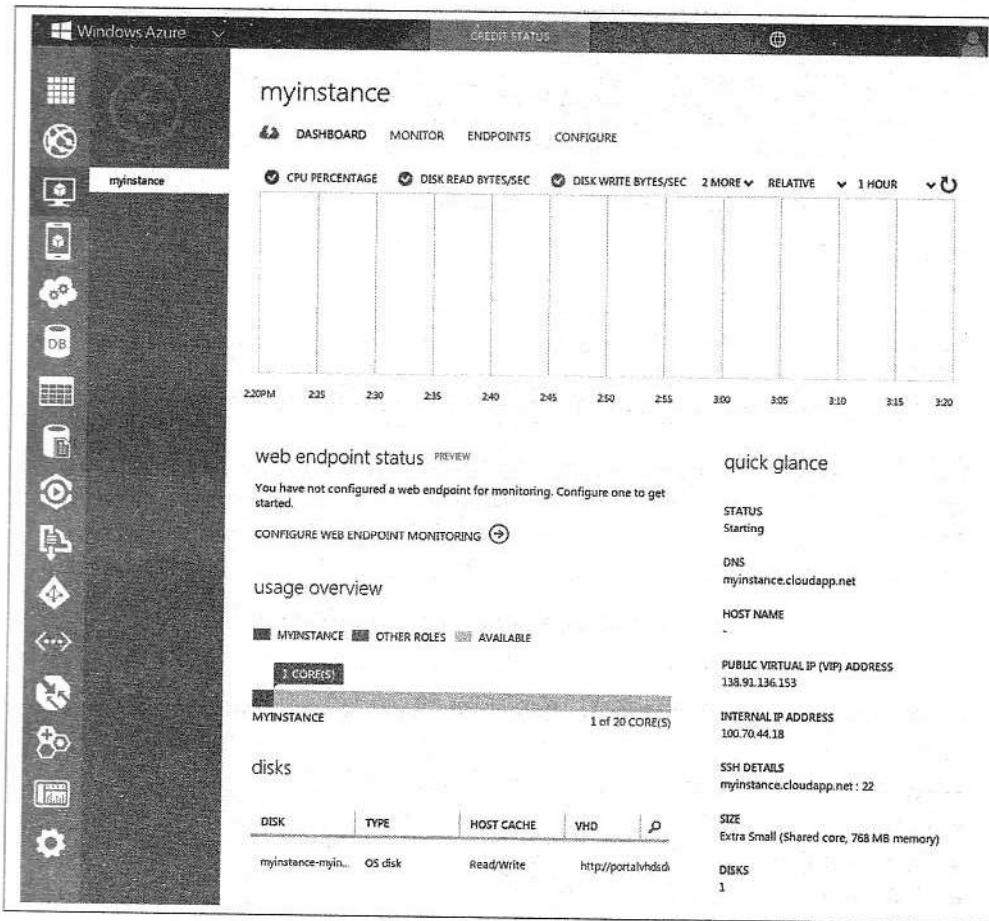


Figure 1.7: Windows Azure Virtual Machines dashboard

smartphones. Sales Cloud allows sales representatives to manage customer profiles, track opportunities, optimize campaigns from lead to close and monitor the impact of campaigns.

Salesforce Service Cloud (TM) is a cloud based customer service management SaaS. Service Cloud provides companies a call-center like view and allows creating, tracking, routing and escalating cases. Service Cloud can be fully integrated with a company's call-center telephony and back office apps. Service Cloud also provides self service capabilities to customers. Service Cloud includes a social networking plug-in that enables social customer service where comments from social media channels can be used to answer customer questions.

Salesforce Marketing Cloud (TM) is cloud based social marketing SaaS. Marketing cloud allows companies to identify sales leads from social media, discover advocates, identify the most trending information on any topic. Marketing Cloud allows companies to pro-actively engage with customers, manage social listening, create and deploy social content, manage and execute optimized social advertisement campaigns and track the performance of social campaigns. Figure 1.9 shows a screenshot of Salesforce dashboard.

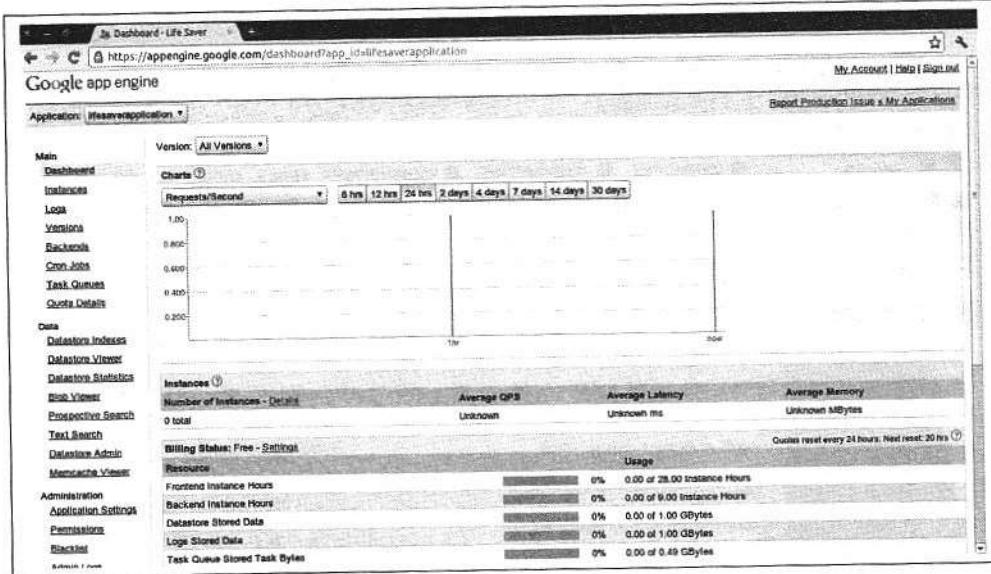


Figure 1.8: Google App Engine dashboard

Some of the tools included in the Salesforce Sales, Service and Marketing Clouds include:

- Accounts and contacts
- Leads
- Opportunities
- Campaigns
- Chatter
- Analytics and Forecasts

1.5 Cloud-based Services & Applications

Having discussed the characteristics, service and deployment models of cloud computing, let us now consider a few examples of the cloud-based services and applications.

1.5.1 Cloud Computing for Healthcare

Figure 1.10 shows the application of cloud computing environments to the healthcare ecosystem [10]. Hospitals and their affiliated providers can securely access patient data stored in the cloud and share the data with other hospitals and physicians. Patients can access their own health information from all of their care providers and store it in a personal health record (PHR) providing them with an integrated record that may even be a family health record. The PHR can be a vehicle for e-prescribing, a technique known to reduce medication dispensing errors and to facilitate medication reconciliation. History and information stored in the cloud (using SaaS applications) can streamline the admissions, care and discharge processes by eliminating redundant data collection and entry. Health payers can increase the effectiveness and lower the cost of their care management programs by providing value added services and giving access to health information to members.



Figure 1.9: Salesforce dashboard

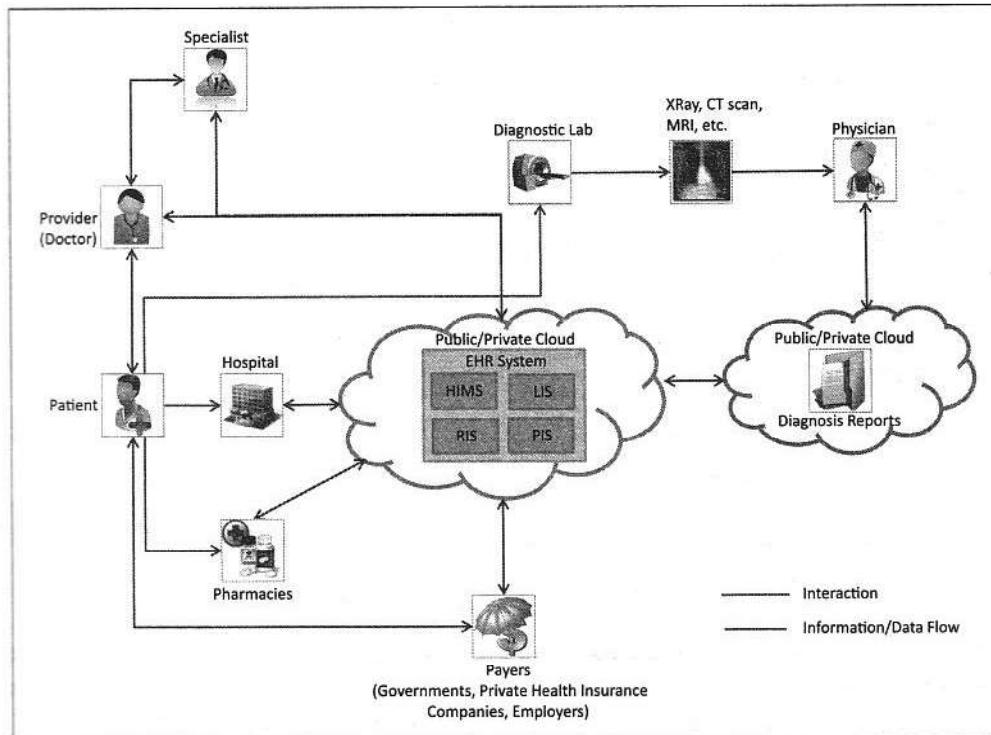


Figure 1.10: Cloud computing for healthcare

1.5.2 Cloud Computing for Energy Systems

Energy systems (such as smart grids, power plants, wind turbine farms, etc.) have thousands of sensors that gather real-time maintenance data continuously for condition monitoring and

failure prediction purposes. These energy systems have a large number of critical components that must function correctly so that the systems can perform their operations correctly. For example, a wind turbine has a number of critical components, e.g., bearings, turning gears, etc. that must be monitored carefully as wear and tear in such critical components or sudden change in operating conditions of the machines can result in failures. In systems such as power grids, real-time information is collected using specialized electrical sensors called Phasor Measurement Units (PMU) at the substations. The information received from PMUs must be monitored in real-time for estimating the state of the system and for predicting failures. Maintenance and repair of such complex systems is not only expensive but also time consuming, therefore failures can cause huge losses for the operators, and supply outage for consumers. In [8], the Bahga & Madisetti have proposed a generic framework, CloudView, for storage, processing and analysis of massive machine maintenance data, collected from a large number of sensors embedded in industrial machines, in a cloud computing environment. The approach proposed in [8], in addition to being the first reported use of the cloud architecture for maintenance data storage, processing and analysis, also evaluated several possible cloud-based architectures that leverage the advantages of the parallel computing capabilities of the cloud to make local decisions with global information efficiently, while avoiding potential data bottlenecks that can occur in getting the maintenance data in and out of the cloud. Figure 1.11 shows a generic use case of cloud for energy systems.

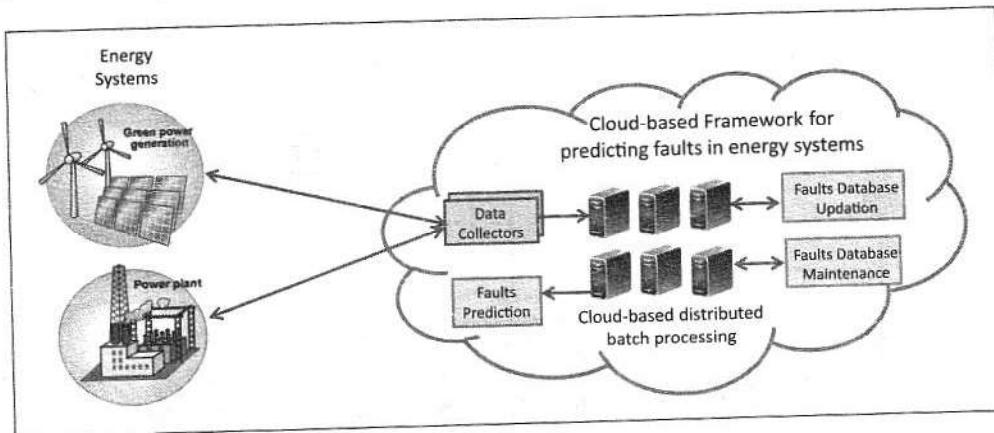


Figure 1.11: Cloud computing for energy systems

1.5.3 Cloud Computing for Transportation Systems

Intelligent transportation systems (ITS) have evolved significantly in recent years. Modern ITS are driven by data collected from multiple sources which is processed to provide new services to their users. By collecting large amount of data from various sources and processing the data into useful information, data-driven ITS can provide new services such as advanced route guidance, dynamic vehicle routing, anticipating customer demands for pickup and delivery problem, etc. Collection and organization of data from multiple sources in real-time and using the massive amounts data for providing intelligent decisions for operations and supply chains, is a major challenge, primarily because the size of the databases involved is very large, and real-time analysis tools have not been available. As a result large organizations

are faced with a seemingly unsurmountable problem of analyzing terabytes of unorganized data stored on isolated and distinct geographical locations. However, recent advances in massive scale data processing systems, utilized for driving business operations of corporations provide a promising approach to massive ITS data storage and analysis.

In recent work, we have proposed a cloud-based framework that can be leveraged for real-time fresh food supply tracking and monitoring [9]. Fresh food can be damaged during transit due to unrefrigerated conditions and changes in environmental conditions such as temperature and humidity, which can lead to microbial infections and biochemical reactions or mechanical damage due to rough handling. Spoilage of fruits and vegetables during transport and distribution not only results in losses to the distributors but also presents a hazard to the food safety. Therefore tracking and monitoring of fresh food supply is an important problem that needs to be addressed. Typically medium and large container trucks are used for fresh food supply.

Since fresh foods have short durability, tracking the supply of fresh foods and monitoring the transit conditions can help identification of potential food safety hazards. The analysis and interpretation of data on the environmental conditions in the container and food truck positioning can enable more effective routing decisions in real time. Therefore, it is possible to take remedial measures such as, (1) the food that has a limited time budget before it gets rotten can be re-routed to a closer destinations, (2) alerts can be raised to the driver and the distributor about the transit conditions, such as container temperature exceeding the allowed limit, humidity levels going out of the allowed limit, etc., and corrective actions can be taken before the food gets damaged. Figure 1.12 shows a generic use case of cloud for transportation systems.

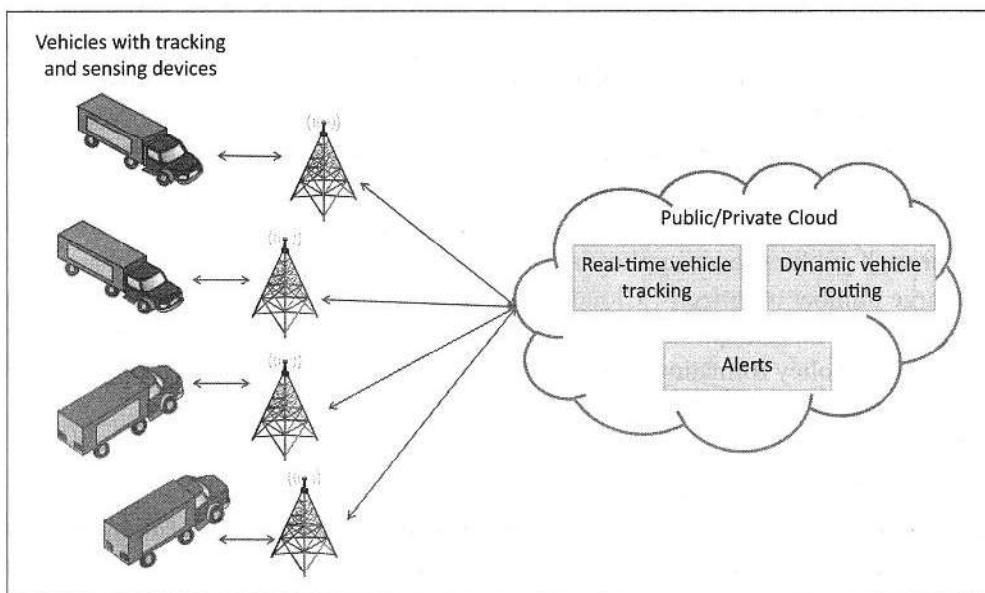


Figure 1.12: Cloud computing for transportation systems

1.5.4 Cloud Computing for Manufacturing Industry

Industrial Control Systems (ICS), such as supervisory control and data acquisition (SCADA) systems, distributed control systems (DCS), and other control system configurations such as Programmable Logic Controllers (PLC) continuously generate monitoring and control data. Real-time collection, management and analysis of data on production operations generated by ICS, in the cloud, can help in estimating the state of the systems, improve plant and personnel safety and thus take appropriate action in real-time to prevent catastrophic failures. Figure 1.13 shows a generic use case of cloud for manufacturing industry.

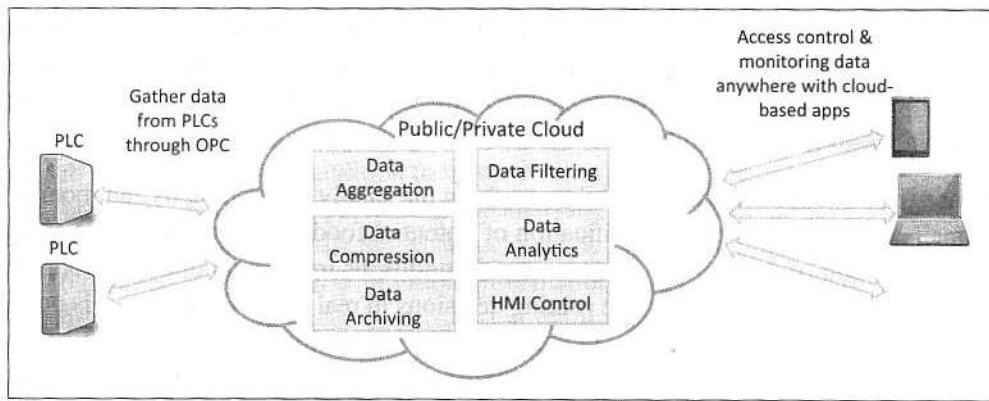


Figure 1.13: Cloud computing for manufacturing industry

1.5.5 Cloud Computing for Government

Cloud computing can play significant role for improving the efficiency and transparency of government operations. Cloud-based e-Governance systems can improve delivery of services to citizens, business, government employees and agencies, etc. and also improve the participation of all responsible parties in various government schemes and policy formation processes. Public services such as public transport reservations, vehicle registrations, issuing of driving licenses, income tax filing, electricity and water bill payments, birth or marriage registration, etc. can be facilitated through cloud-based applications. The benefit of using cloud for such public service applications is that the applications can be scaled up to serve a very large number of citizens. Cloud-based applications can share common data related to citizens. Data on utilization of government schemes can be collected from the citizens and used in the policy formation process and improvement of schemes. Figure 1.14 shows a generic use case of cloud for government.

1.5.6 Cloud Computing for Education

Cloud computing can help in improving the reach of quality online education to students. Cloud-based collaboration applications such as online forums, can help student discuss common problems and seek guidance from experts. Universities, colleges and schools can use cloud-based information management systems to admissions, improve administrative efficiency, offer online and distance education programs, online exams, track progress of students, collect feedback from students, for instance. Cloud-based online learning systems

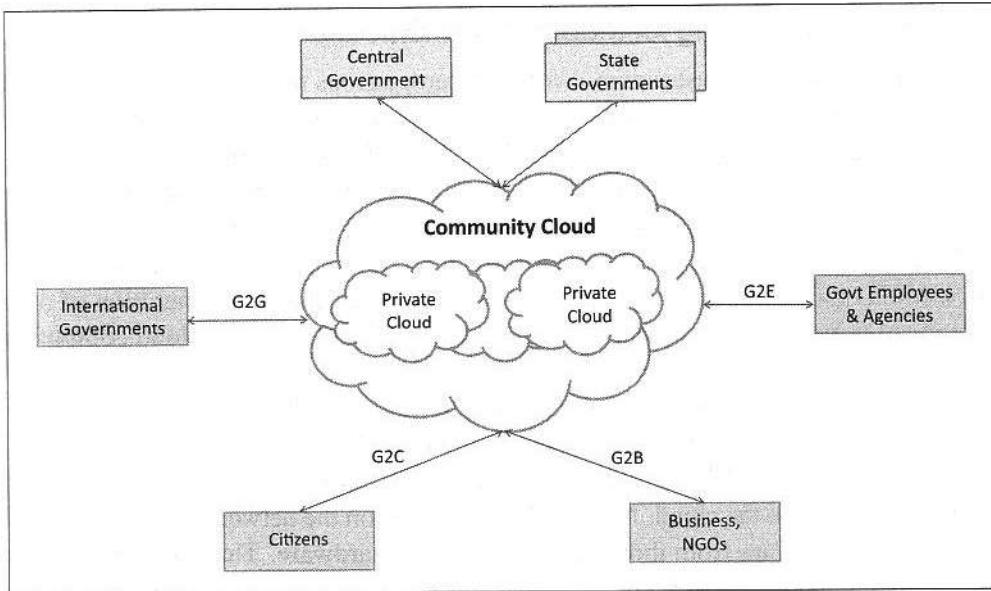


Figure 1.14: Cloud computing for government

can provide access to high quality educational material to students. Figure 1.15 shows a generic use case of cloud for education. Cloud-based systems can help universities, colleges and schools in cutting down the IT infrastructure costs and yet provide access to educational services to a large number of students.

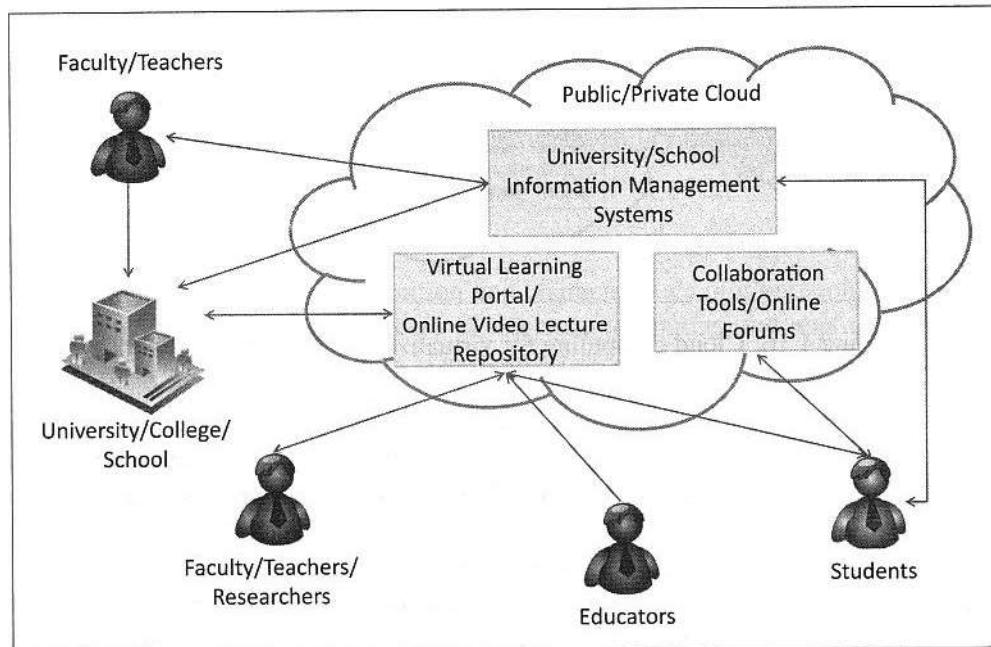


Figure 1.15: Cloud computing for education

1.5.7 Cloud Computing for Mobile Communication

Mobile communication infrastructure involves heterogeneous network devices for the radio access network (RAN) and the core network (CN). A variety of proprietary hardware components and systems are used for these network devices adding to their cost and inflexibility. Expansion and upgradation of the mobile network requires significant capital investments to meet the hardware and space requirements. Due to the increasing speed of innovation, the lifecycles of the network devices are becoming shorter. Network Function Virtualization (NFV) is being seen as a key enabling technology for the fifth generation of mobile communication networks (5G) in the next decade. NFV will leverage cloud computing to consolidate the heterogeneous network devices into the cloud. The NFV architecture, as being standardized by the European Telecommunications Standards Institute (ETSI) comprises of NFV infrastructure, virtual network functions and NFV management and orchestration layers [11]. NFV comprises of network functions implemented in software that run on virtualized resources in the cloud. NFV enables a separation the network functions which are implemented in software from those of the underlying hardware. Thus, network functions can be easily tested and upgraded by installing new software while the hardware remains the same. This flexibility will speed up innovation and reduce the time-to-market. By leveraging the cloud for mobile communication network functions significant savings in capital and operational expenditure can be achieved.

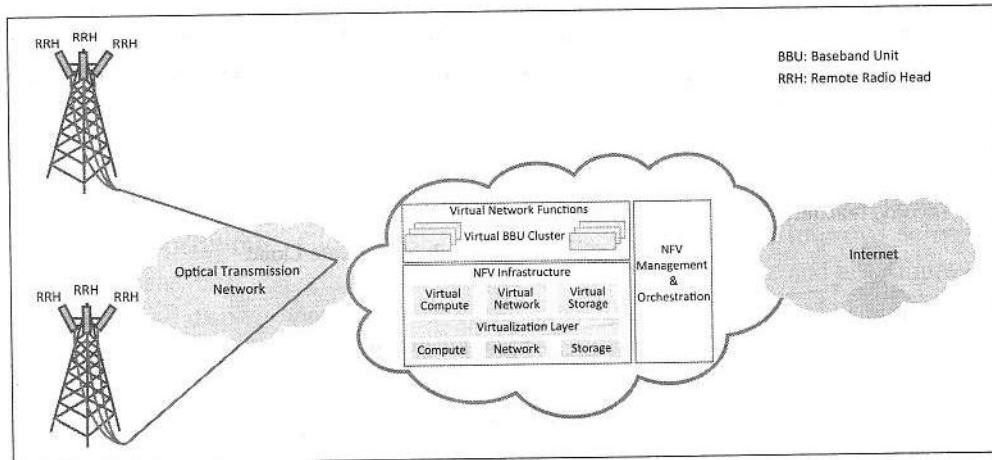


Figure 1.16: Cloud computing for virtualizing radio access network

Figure 1.16 shows a use case of cloud-based NFV architecture for cloud-based radio access networks (C-RANs) with virtualized mobile base stations (baseband units). The baseband units (BBUs), such as eNodeB in 4G, in current mobile communication networks are co-located with the cell towers on-site and run on proprietary hardware. The BBUs are typically designed for worst-case peak loads. However, typical workload levels are much lower than the peak loads, therefore, the excess capacity goes unused. With NFV and cloud the BBUs can be virtualized and only as many resources as required to meet the workload levels can be provisioned on-demand. This will result in significant power savings. Centralized cloud-based virtual BBU clusters can replace on-site installations of BBUs in distributed geographical locations. This will result in reduction of management and

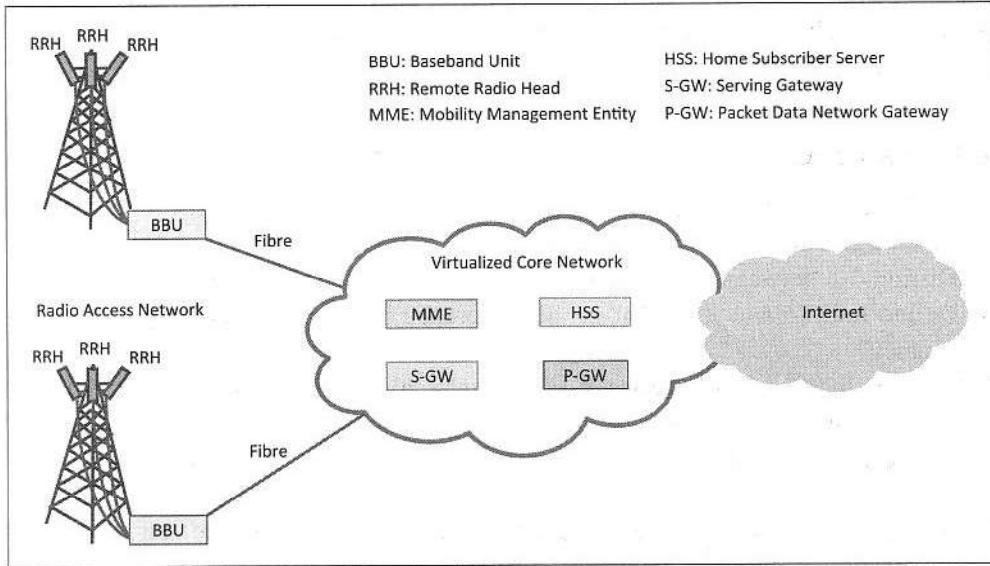


Figure 1.17: Cloud computing for virtualizing mobile core network

operational expenses.

Figure 1.17 shows a use case of cloud-based NFV architecture for mobile core network. With NFV, the core network devices such as Mobility Management Entity (MME), Home Subscriber Server (HSS), Serving Gateway (S-GW) and Packet Data Network Gateway (P-GW) in 4G can be implemented in software and deployed on virtualized resources in the cloud. This will reduce the total cost of ownership due to consolidation of network component that run on industry standard networking hardware. Other benefits of using cloud-based NFV architecture for mobile core network include improved resource utilization efficiency, improved network resilience, improved flexibility in scaling up capacity.

Summary

In this chapter you learned the definition and characteristics of cloud computing. Cloud computing offers Internet-based access to low cost computing and applications that are provided using virtualized resources. On-demand service, remote accessibility through a variety of networks, resource pooling, rapid elasticity and measured service are the key characteristics of cloud computing. Cloud computing resources can be provisioned on-demand by the users. Cloud computing resources can be accessed over the network with standard access mechanisms. Cloud resources are pooled to serve multiple users using multi-tenancy.

Cloud computing has three service models - IaaS, PaaS and SaaS. IaaS provides the users the capability to provision computing and storage resources. PaaS provides the users the capability to develop and deploy their own applications in the cloud. SaaS provides applications hosted in the cloud through thin client interfaces.

Cloud computing is being increasingly adopted by individual users, small and large enterprises, large organizations and governments. Cloud computing is being applied in

various fields such as healthcare, education, governance, energy systems, manufacturing industry, transportation systems, etc.

Review Questions

1. Define cloud computing
2. List the pros and cons of cloud computing.
3. Distinguish between IaaS, PaaS and SaaS.
4. Define multi-tenancy. What is the difference between virtual and organic multi-tenancy?
5. What is the difference between horizontal scaling and vertical scaling? Describe scenarios in which you will use each type of scaling.
6. Define virtualization. What is the difference between full, para- and hardware-assisted virtualization?
7. Assume your company wants to launch an e-commerce website. Which cloud services and deployment models will you consider for the website?

2 — Cloud Concepts & Technologies

This Chapter Covers

Concepts and enabling technologies of cloud computing including:

- Virtualization
- Load balancing
- Scalability & Elasticity
- Deployment
- Replication
- Monitoring
- MapReduce
- Identity and Access Management
- Service Level Agreements
- Billing

In this chapter you will learn the key concepts and enabling technologies of cloud computing. We will introduce and build upon technologies such as virtualization, load balancing, and on-demand provisioning. A popular programming model, called MapReduce, will also be covered.

2.1 Virtualization

Virtualization refers to the partitioning of the resources of a physical system (such as computing, storage, network and memory) into multiple virtual resources. Virtualization is the key enabling technology of cloud computing and allows pooling of resources. In cloud computing, resources are pooled to serve multiple users using multi-tenancy. Multi-tenant aspects of the cloud allow multiple users to be served by the same physical hardware. Users are assigned virtual resources that run on top of the physical resources. Figure 2.1 shows the architecture of a virtualization technology in cloud computing. The physical resources such as computing, storage memory and network resources are virtualized. The virtualization layer partitions the physical resources into multiple virtual machines. The virtualization layer allows multiple operating system instances to run currently as virtual machines on the same underlying physical resources.

Hypervisor

The virtualization layer consists of a hypervisor or a virtual machine monitor (VMM). The hypervisor presents a virtual operating platform to a guest operating system (OS). There are two types of hypervisors as shown in Figures 2.2 and 2.3. Type-1 hypervisors or the native hypervisors run directly on the host hardware and control the hardware and monitor the guest operating systems. Type 2 hypervisors or hosted hypervisors run on top of a conventional (main/host) operating system and monitor the guest operating systems.

Guest OS

A guest OS is an operating system that is installed in a virtual machine in addition to the host or main OS. In virtualization, the guest OS can be different from the host OS.

Various forms of virtualization approaches exist:

Full Virtualization

In full virtualization, the virtualization layer completely decouples the guest OS from the underlying hardware. The guest OS requires no modification and is not aware that it is being virtualized. Full virtualization is enabled by direct execution of user requests and binary translation of OS requests. Figure 2.4 shows the full virtualization approach.

Para-Virtualization

In para-virtualization, the guest OS is modified to enable communication with the hypervisor to improve performance and efficiency. The guest OS kernel is modified to replace non-virtualizable instructions with hypercalls that communicate directly with the virtualization layer hypervisor. Figure 2.5 shows the para-virtualization approach.

Hardware Virtualization

Hardware assisted virtualization is enabled by hardware features such as Intel's Virtualization Technology (VT-x) and AMD's AMD-V. In hardware assisted virtualization, privileged and

sensitive calls are set to automatically trap to the hypervisor. Thus, there is no need for either binary translation or para-virtualization.

Table 2.1 lists some examples of popular hypervisors.

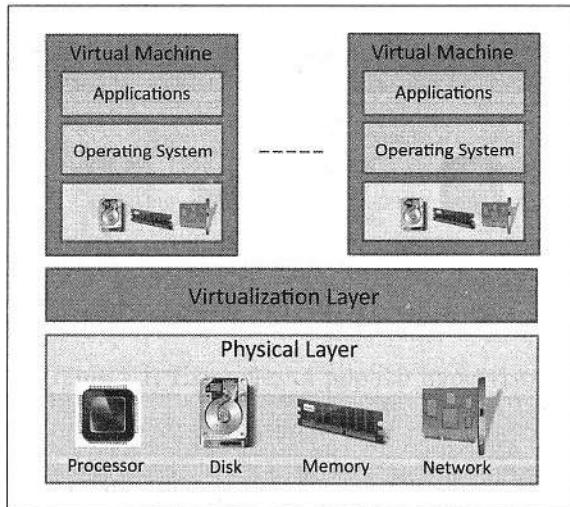


Figure 2.1: Virtualization architecture

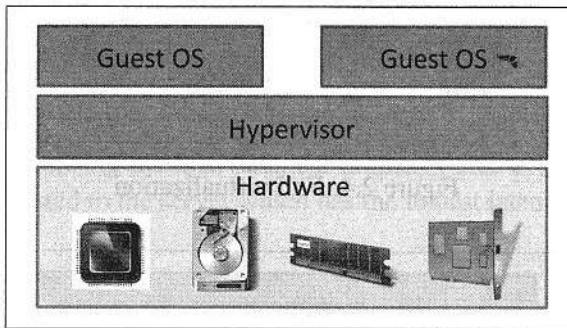


Figure 2.2: Hypervisor design: Type-1

2.2 Load Balancing

One of the important features of cloud computing is scalability. Cloud computing resources can be scaled up on demand to meet the performance requirements of applications. Load balancing distributes workloads across multiple servers to meet the application workloads. The goals of load balancing techniques are to achieve maximum utilization of resources, minimizing the response times, maximizing throughput. Load balancing distributes the incoming user requests across multiple resources. With load balancing, cloud-based applications can achieve high availability and reliability. Since multiple resources under a load balancer are used to serve the user requests, in the event of failure of one or more of the resources, the load balancer can automatically reroute the user traffic to the healthy resources. To the end user accessing a cloud-based application, a load balancer makes the pool of servers under the

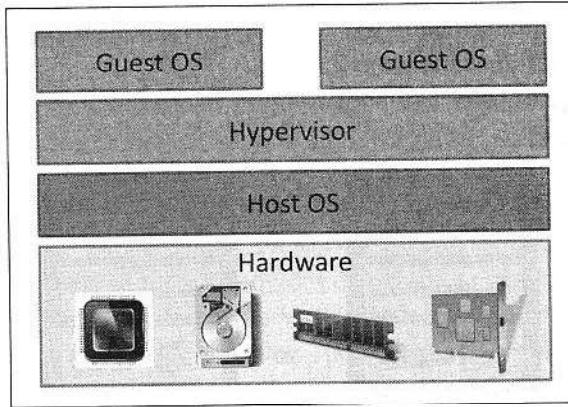


Figure 2.3: Hypervisor design: Type-2

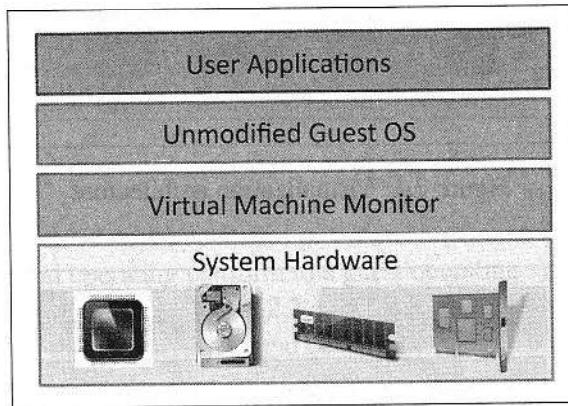


Figure 2.4: Full virtualization

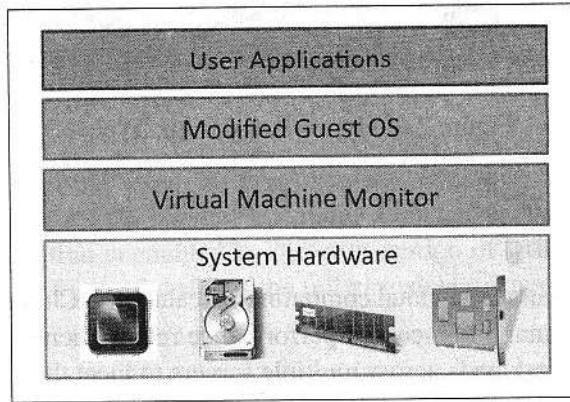


Figure 2.5: Para-virtualization

load balancer appear as a single server with high computing capacity. The routing of user requests is determined based on a load balancing algorithm. Commonly used load balancing algorithms include:

Hypervisor	Type
Citrix XenServer	Type-1
Oracle VM Server	Type-1
KVM	Type-1
VMWare ESX/ESXi	Type-1
Microsoft Hyper-V	Type-1
Xen Hypervisor	Type-1
VMWare Workstation	Type-2
VirtualBox	Type-2

Table 2.1: Examples of popular hypervisors

Round Robin

In round robin load balancing, the servers are selected one by one to serve the incoming requests in a non-hierarchical circular fashion with no priority assigned to a specific server.

Weighted Round Robin

In weighted round robin load balancing, servers are assigned some weights. The incoming requests are proportionally routed using a static or dynamic ratio of respective weights.

Low Latency

In low latency load balancing the load balancer monitors the latency of each server. Each incoming request is routed to the server which has the lowest latency.

Least Connections

In least connections load balancing, the incoming requests are routed to the server with the least number of connections.

Priority

In priority load balancing, each server is assigned a priority. The incoming traffic is routed to the highest priority server as long as the server is available. When the highest priority server fails, the incoming traffic is routed to a server with a lower priority.

Overflow

Overflow load balancing is similar to priority load balancing. When the incoming requests to highest priority server overflow, the requests are routed to a lower priority server.

Figure 2.6 depicts these various load balancing approaches. For session based applications, an important issue to handle during load balancing is the persistence of multiple requests from a particular user session. Since load balancing can route successive requests from a user session to different servers, maintaining the state or the information of the session is important. Three commonly used persistence approaches are described below:

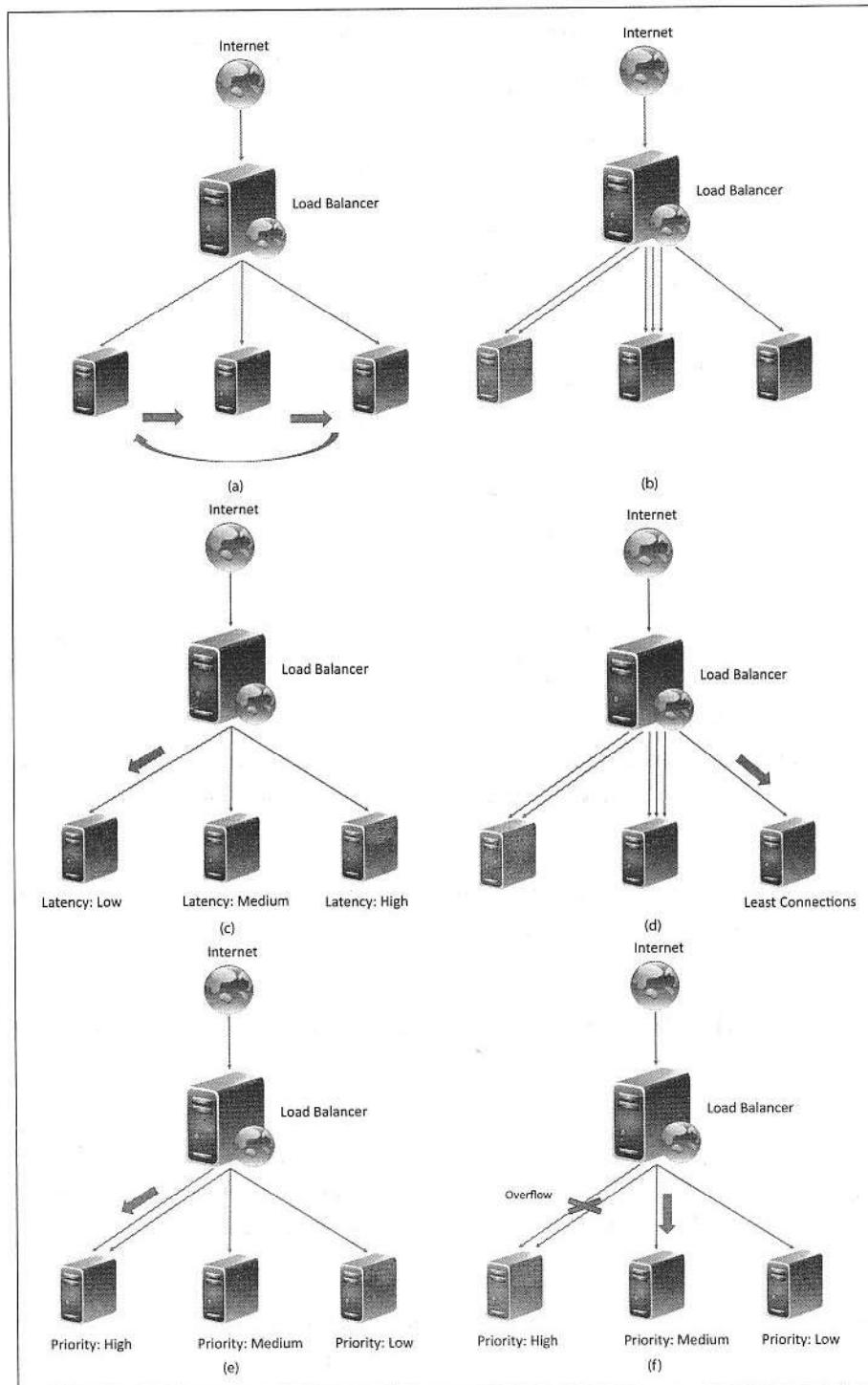


Figure 2.6: (a) Round robin load balancing, (b) Weighted round robin load balancing, (c) Low latency load balancing, (d) Least connections load balancing, (e) Priority load balancing, (f) Overload load balancing

Sticky sessions

In this approach all the requests belonging to a user session are routed to the same server. These sessions are called sticky sessions. The benefit of this approach is that it makes session management simple. However, a drawback of this approach is that if a server fails all the sessions belonging to that server are lost, since there is no automatic failover possible.

Session Database

In this approach, all the session information is stored externally in a separate session database, which is often replicated to avoid a single point of failure. Though, this approach involves additional overhead of storing the session information, however, unlike the sticky session approach, this approach allows automatic failover.

Browser cookies

In this approach, the session information is stored on the client side in the form of browser cookies. The benefit of this approach is that it makes the session management easy and has the least amount of overhead for the load balancer.

URL re-writing

In this approach, a URL re-write engine stores the session information by modifying the URLs on the client side. Though this approach avoids overhead on the load balancer, a drawback is that the amount of session information that can be stored is limited. For applications that require larger amounts of session information, this approach does not work.

Load balancing can be implemented in software or hardware. Software-based load balancers run on standard operating systems, and like other cloud resources, load balancers are also virtualized. Hardware-based load balancers implement load balancing algorithms in Application Specific Integrated Circuits (ASICs). In a hardware load balancer, the incoming user requests are routed to the underlying servers based on some pre-configured load balancing strategy and the response from the servers are sent back either directly to the user (at layer-4) or back to the load balancer (at layer-7) where it is manipulated before being sent back to the user. Table 2.2 lists some examples of load balancers.

2.3 Scalability & Elasticity

Multi-tier applications such as e-Commerce, social networking, business-to-business, etc. can experience rapid changes in their traffic. Each website has a different traffic pattern which is determined by a number of factors that are generally hard to predict beforehand. Modern web applications have multiple tiers of deployment with varying number of servers in each tier. Capacity planning is an important task for such applications. Capacity planning involves determining the right sizing of each tier of the deployment of an application in terms of the number of resources and the capacity of each resource. Capacity planning may be for computing, storage, memory or network resources. Figure 2.7 shows the cost versus capacity curves for traditional and cloud approaches.

Traditional approaches for capacity planning are based on predicted demands for applications and account for worst case peak loads of applications. When the workloads of applications increase, the traditional approaches have been either to scale up or scale

Load Balancer	Type
Nginx	Software
HAProxy	Software
Pound	Software
Varish	Software
Cisco Systems Catalyst 6500	Hardware
Coyote Point Equalizer	Hardware
F5 Networks BIG-IP LTM	Hardware
Barracuda Load Balancer	Hardware

Table 2.2: Examples of popular load balancers

out. Scaling up involves upgrading the hardware resources (adding additional computing, memory, storage or network resources). Scaling out involves addition of more resources of the same type. Traditional scaling up and scaling out approaches are based on demand forecasts at regular intervals of time. When variations in workloads are rapid, traditional approaches are unable to keep track with the demand and lead to either over-provisioning or under-provisioning of resources. Over-provisioning of resources leads to higher capital expenditures than required. On the other hand, under-provisioning of resources leads to traffic overloads, slow response times, low throughputs and hence loss of opportunity to serve the customers. Analyzing the real traffic history plots for top websites shown in Figure 2.7 we observe that the off peak workloads are significantly lower than peak workloads. Traditional capacity planning approaches which are designed to meet the peak loads result in excess capacity and under utilization of resources. Moreover, the infrastructure resources for traditional applications are fixed, rigid and provisioned in advance. This involves up-front capital expenditures for setting up the infrastructure.

2.4 Deployment

Figure 2.8 shows the cloud application deployment lifecycle. Deployment prototyping can help in making deployment architecture design choices. By comparing performance of alternative deployment architectures, deployment prototyping can help in choosing the best and most cost effective deployment architecture that can meet the application performance requirements. Table 2.3 lists some popular cloud deployment management tools. Deployment design is an iterative process that involves the following steps:

Deployment Design

In this step the application deployment is created with various tiers as specified in the deployment configuration. The variables in this step include the number of servers in each tier, computing, memory and storage capacities of servers, server interconnection, load balancing and replication strategies. Deployment is created by provisioning the cloud

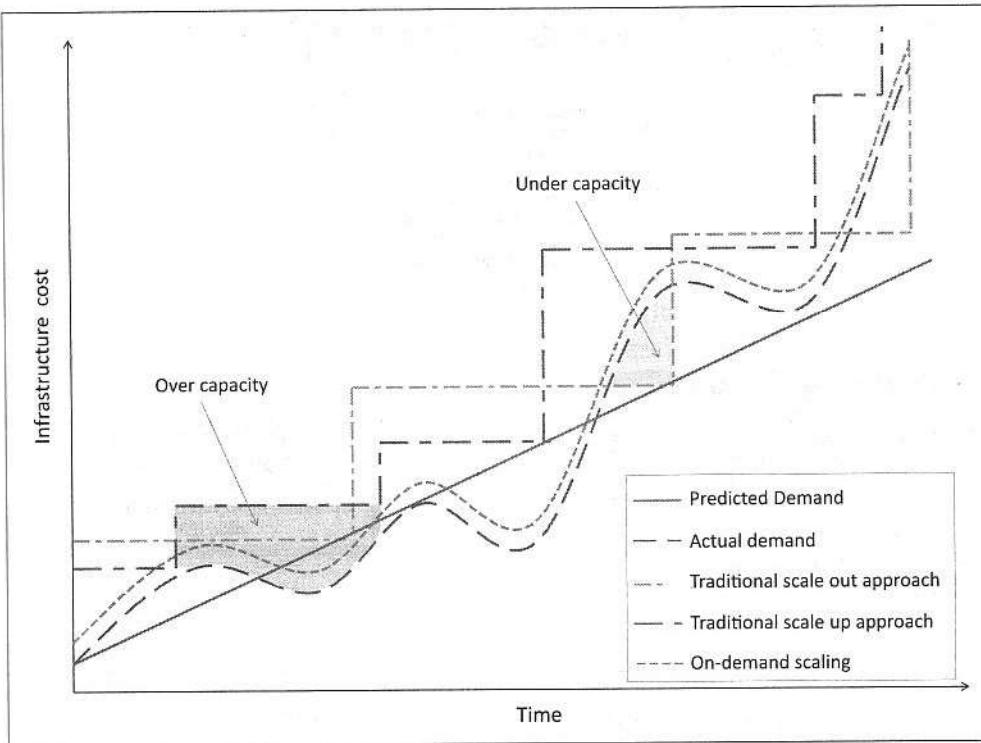


Figure 2.7: Cost versus capacity curves

resources as specified in the deployment configuration. The process of resource provisioning and deployment creation is often automated and involves a number of steps such as launching of server instances, configuration of servers, and deployment of various tiers of the application on the servers.

Performance Evaluation

Once the application is deployed in the cloud, the next step in the deployment lifecycle is to verify whether the application meets the performance requirements with the deployment. This step involves monitoring the workload on the application and measuring various workload parameters such as response time and throughput. In addition to this, the utilization of servers (CPU, memory, disk, I/O, etc.) in each tier is also monitored.

Deployment Refinement

After evaluating the performance of the application, deployments are refined so that the application can meet the performance requirements. Various alternatives can exist in this step such as vertical scaling (or scaling up), horizontal scaling (or scaling out), alternative server interconnections, alternative load balancing and replication strategies, for instance.

2.5 Replication

Replication is used to create and maintain multiple copies of the data in the cloud. Replication of data is important for practical reasons such as business continuity and disaster recovery.

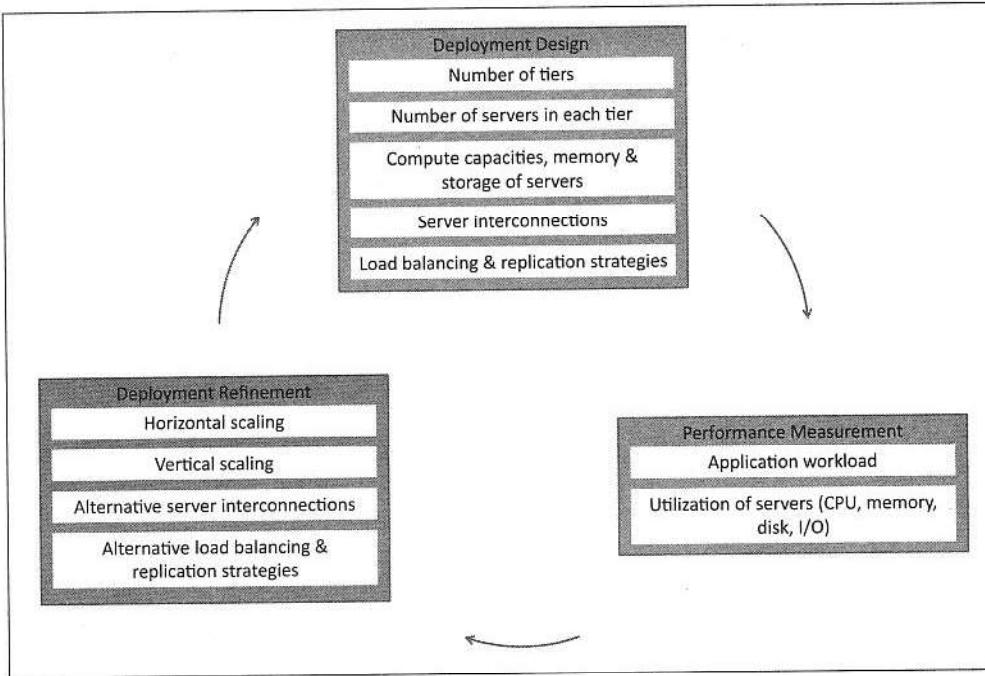


Figure 2.8: Cloud application deployment lifecycle

In the event of data loss at the primary location, organizations can continue to operate their applications from secondary data sources. With real-time replication of data, organizations can achieve faster recovery from failures. Traditional business continuity and disaster recovery approaches don't provide efficient, cost effective and automated recovery of data. Cloud based data replication approaches provide replication of data in multiple locations, automated recovery, low recovery point objective (RPO) and low recovery time objective (RTO). Cloud enables rapid implementation of replication solutions for disaster recovery for small and medium enterprises and large organizations. With cloud-based data replication organizations can plan for disaster recovery without making any capital expenditures on purchasing, configuring and managing secondary site locations. Cloud provides affordable replication solutions with pay-per-use/pay-as-you-go pricing models. There are three types of replication approaches as shown in Figure 2.9 and described as follows:

Array-based Replication

Array-based replication uses compatible storage arrays to automatically copy data from a local storage array to a remote storage array. Arrays replicate data at the disk sub-system level, therefore the type of hosts accessing the data and the type of data is not important. Thus array-based replication can work in heterogeneous environments with different operating systems. Array-based replication uses Network Attached Storage (NAS) or Storage Area Network (SAN), to replicate. A drawback of this array-based replication is that it requires similar arrays at local and remote locations. Thus the costs for setting up array-based replication are higher than the other approaches.

Cloud Deployment Management Tool	Features
RightScale	Design, deploy and manage cloud deployments across multiple public or private clouds.
Scalr	Provides tools to automate the management of servers, monitors servers, replaces servers that fail, provides auto scaling and backups.
Kaavo	Allows deploying applications easily across multiple clouds, managing distributed applications and automating high availability.
CloudStack	Allows simple and cost effective deployment management and configuration of cloud computing environments.

Table 2.3: Examples of popular cloud deployment management tools

Network-based Replication

Network-based replication uses an appliance that sits on the network and intercepts packets that are sent from hosts and storage arrays. The intercepted packets are replicated to a secondary location. The benefits of this approach is that it supports heterogeneous environments and requires a single point of management. However, this approach involves higher initial costs due to replication hardware and software.

Host-based Replication

Host-based replication runs on standard servers and uses software to transfer data from a local to remote location. The host acts the replication control mechanism. An agent is installed on the hosts that communicates with the agents on the other hosts. Host-based replication can either be block-based or file-based. Block-based replication typically require dedicated volumes of the same size on both the local and remote servers. File-based replication requires less storage as compared to block-based storage. File-based replication gives additional allows the administrators to choose the files or folders to be replicated. Host-based replication with cloud-infrastructure provides affordable replication solutions. With host-based replication, entire virtual machines can be replicated in real-time.

2.6 Monitoring

Cloud resources can be monitored by monitoring services provided by the cloud service providers. Monitoring services allow cloud users to collect and analyze the data on various monitoring metrics. Figure 2.10 shows a generic architecture for a cloud monitoring service. A monitoring service collects data on various system and application metrics from the cloud computing instances. Monitoring services provide various pre-defined metrics. Users can also define their custom metrics for monitoring the cloud resources. Users can define various actions based on the monitoring data, for example, auto-scaling a cloud deployment when the CPU usage of monitored resources becomes high. Monitoring services also provide various statistics based on the monitoring data collected. Table 2.4 lists the commonly

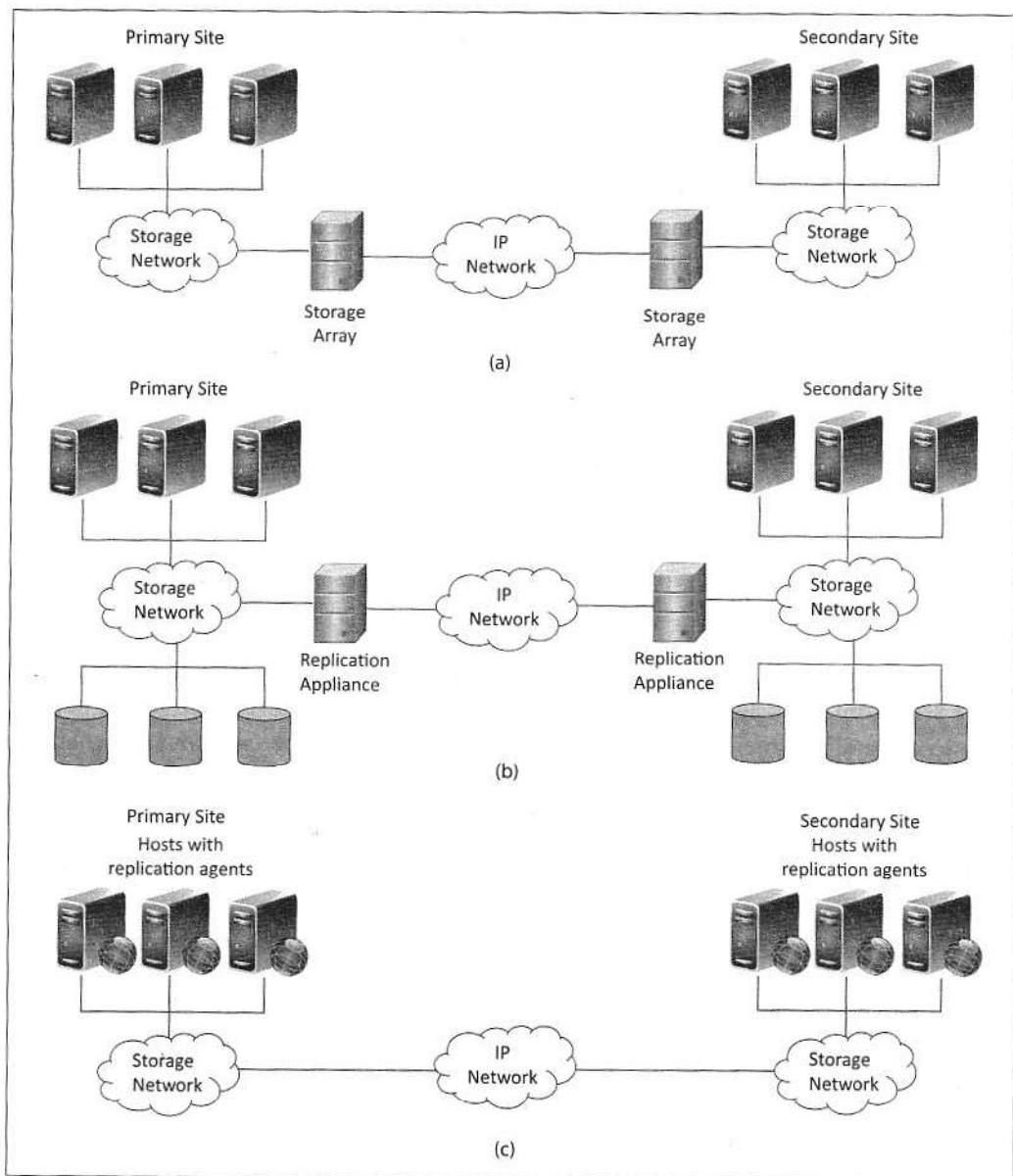


Figure 2.9: Replication approaches: (a) Array-based replication, (b) Network-based replication, (c) Host-based replication

used monitoring metrics for cloud computing resources. Monitoring of cloud resources is important because it allows the users to keep track of the health of applications and services deployed in the cloud. For example, an organization which has its website hosted in the cloud can monitor the performance of the website and also the website traffic. With the monitoring data available at run-time users can make operational decisions such as scaling up or scaling down cloud resources.

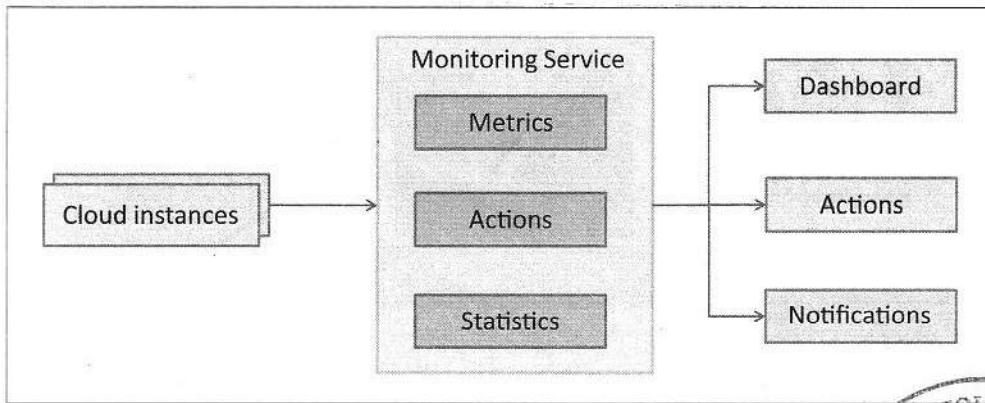
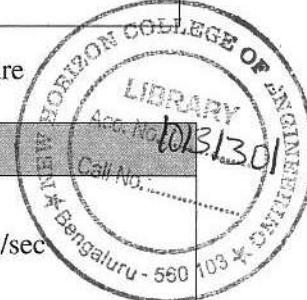


Figure 2.10: Typical cloud monitoring service architecture



Type	Metrics
CPU	CPU-Usage, CPU-Idle
Disk	Disk-Usage, Bytes/sec (read/write), Operations/sec
Memory	Memory-Used, Memory-Free, Page-Cache
Interface	Packets/sec (incoming/outgoing), Octets/sec(incoming/outgoing)

Table 2.4: Typical monitoring metrics

2.7 Software Defined Networking

Software-Defined Networking (SDN) is a networking architecture that separates the control plane from the data plane and centralizes the network controller. Figure 2.11 shows the conventional network architecture built with specialized hardware (switches, routers, etc.). Network devices in conventional network architectures are getting exceedingly complex with the increasing number of distributed protocols being implemented and the use of proprietary hardware and interfaces. In the conventional network architecture the control plane and data plane are coupled. Control plane is the part of the network that carries the signaling and routing message traffic while the data plane is the part of the network that carries the payload data traffic.

The limitations of the conventional network architectures are as follows:

- **Complex Network Devices:** Conventional networks are getting increasingly complex with more and more protocols being implemented to improve link speeds and reliability. Interoperability is limited due to the lack of standard and open interfaces. Network devices use proprietary hardware and software and have slow product lifecycles limiting innovation. The conventional networks were well suited for static traffic patterns and had a large number of protocols designed for specific applications. With the emergence of cloud computing and proliferation of internet access devices, the traffic patterns are becoming more and more dynamic. Due to the complexity of conventional network devices, making changes in the networks to meet the dynamic traffic patterns has