



## Cloud computing and emerging IT platforms: Vision, hype, and reality for delivering computing as the 5th utility

Rajkumar Buyya<sup>a,b,\*</sup>, Chee Shin Yeo<sup>a</sup>, Srikumar Venugopal<sup>a</sup>, James Broberg<sup>a</sup>, Ivona Brandic<sup>c</sup>

<sup>a</sup> Grid Computing and Distributed Systems (GRIDS) Laboratory, Department of Computer Science and Software Engineering, The University of Melbourne, Australia

<sup>b</sup> Manjrasoft Pty Ltd, Melbourne, Australia

<sup>c</sup> Institute of Information Systems, Vienna University of Technology, Argentinierstraße 8, 1040 Vienna, Austria

### ARTICLE INFO

#### Article history:

Received 23 September 2008

Received in revised form

21 November 2008

Accepted 3 December 2008

Available online 11 December 2008

#### Keywords:

Cloud computing

Data Centers

Utility computing

Virtualization

Market-oriented resource allocation

### ABSTRACT

With the significant advances in Information and Communications Technology (ICT) over the last half century, there is an increasingly perceived vision that computing will one day be the 5th utility (after water, electricity, gas, and telephony). This computing utility, like all other four existing utilities, will provide the basic level of computing service that is considered essential to meet the everyday needs of the general community. To deliver this vision, a number of computing paradigms have been proposed, of which the latest one is known as Cloud computing. Hence, in this paper, we define Cloud computing and provide the architecture for creating Clouds with market-oriented resource allocation by leveraging technologies such as Virtual Machines (VMs). We also provide insights on market-based resource management strategies that encompass both customer-driven service management and computational risk management to sustain Service Level Agreement (SLA)-oriented resource allocation. In addition, we reveal our early thoughts on interconnecting Clouds for dynamically creating global Cloud exchanges and markets. Then, we present some representative Cloud platforms, especially those developed in industries, along with our current work towards realizing market-oriented resource allocation of Clouds as realized in Aneka enterprise Cloud technology. Furthermore, we highlight the difference between High Performance Computing (HPC) workload and Internet-based services workload. We also describe a meta-negotiation infrastructure to establish global Cloud exchanges and markets, and illustrate a case study of harnessing 'Storage Clouds' for high performance content delivery. Finally, we conclude with the need for convergence of competing IT paradigms to deliver our 21st century vision.

© 2008 Elsevier B.V. All rights reserved.

### 1. Introduction

Computing is being transformed to a model consisting of services that are commoditized and delivered in a manner similar to traditional utilities such as water, electricity, gas, and telephony. In such a model, users access services based on their requirements without regard to where the services are hosted or how they are delivered. Several computing paradigms have promised to deliver this *utility computing* vision and these include cluster computing, Grid computing, and more recently *Cloud computing*. The latter term denotes the infrastructure as a “*Cloud*” from which

businesses and users are able to access applications from anywhere in the world on demand. Thus, the computing world is rapidly transforming towards developing software for millions to consume as a service, rather than to run on their individual computers.

At present, it is common to access content across the Internet independently without reference to the underlying hosting infrastructure. This infrastructure consists of data centers that are monitored and maintained around the clock by content providers. Cloud computing is an extension of this paradigm wherein the capabilities of business applications are exposed as sophisticated services that can be accessed over a network. Cloud service providers are incentivized by the profits to be made by charging consumers for accessing these services. Consumers, such as enterprises, are attracted by the opportunity for reducing or eliminating costs associated with “in-house” provision of these services. However, since cloud applications may be crucial to the core business operations of the consumers, it is essential that the consumers have guarantees from providers on service delivery. Typically, these are provided through Service Level Agreements (SLAs) brokered between the providers and consumers.

\* Corresponding address: Grid Computing and Distributed Systems (GRIDS) Laboratory, Department of Computer Science and Software Engineering, The University of Melbourne, ICT Building, 111, Barry Street, VIC 3053 Melbourne, Victoria, Australia. Tel.: +61 3 83441344.

E-mail addresses: [raj@csse.unimelb.edu.au](mailto:raj@csse.unimelb.edu.au) (R. Buyya), [csyeo@csse.unimelb.edu.au](mailto:csyeo@csse.unimelb.edu.au) (C.S. Yeo), [srikumar@csse.unimelb.edu.au](mailto:srikumar@csse.unimelb.edu.au) (S. Venugopal), [brobergj@csse.unimelb.edu.au](mailto:brobergj@csse.unimelb.edu.au) (J. Broberg), [ivona@infosys.tuwien.ac.at](mailto:ivona@infosys.tuwien.ac.at) (I. Brandic).

Providers such as Amazon, Google, Salesforce, IBM, Microsoft, and Sun Microsystems have begun to establish new data centers for hosting Cloud computing applications in various locations around the world to provide redundancy and ensure reliability in case of site failures. Since user requirements for cloud services are varied, service providers have to ensure that they can be flexible in their service delivery while keeping the users isolated from the underlying infrastructure. Recent advances in microprocessor technology and software have led to the increasing ability of commodity hardware to run applications within *Virtual Machines* (VMs) efficiently. VMs allow both the isolation of applications from the underlying hardware and other VMs, and the customization of the platform to suit the needs of the end-user. Providers can expose applications running within VMs, or provide access to VMs themselves as a service (e.g. Amazon Elastic Compute Cloud) thereby allowing consumers to install their own applications. While convenient, the use of VMs gives rise to further challenges such as the intelligent allocation of physical resources for managing competing resource demands of the users.

In addition, enterprise service consumers with global operations require faster response time, and thus save time by distributing workload requests to multiple Clouds in various locations at the same time. This creates the need for establishing a computing atmosphere for dynamically interconnecting and provisioning Clouds from multiple domains within and across enterprises. There are many challenges involved in creating such Clouds and Cloud interconnections.

Therefore, this paper discusses the current trends in the space of Cloud computing and presents candidates for future enhancements of this technology. This paper is primarily divided into two parts. The first part examines current research issues and developments by:

- presenting the 21st century vision of computing and describing various computing paradigms that have promised or are promising to deliver this grand vision (Section 2),
- differentiating Cloud computing from two other widely explored computing paradigms: Cluster computing and Grid computing (Section 3),
- focusing on VM-centric Cloud services and presenting an architecture for creating market-oriented Clouds using VMs (Section 4),
- providing insights on market-based resource management strategies that encompass both customer-driven service management and computational risk management to sustain SLA-oriented resource allocation (Section 5),
- revealing our early thoughts on interconnecting Clouds for dynamically creating global Cloud exchanges and markets (Section 6), and
- comparing some representative Cloud platforms, especially those developed in industries along with our Aneka enterprise Cloud technology (Section 7).

The second part introduces our current work on Cloud computing which include:

- realizing market-oriented resource allocation of Clouds as realized in Aneka enterprise Cloud technology and highlighting the difference between High Performance Computing (HPC) workload and Internet-based services workload (Section 8),
- incorporating a meta-negotiation infrastructure for QoS management to establish global Cloud exchanges and markets (Section 9), and
- creating 3rd party cloud services based on high performance content delivery over commercial cloud storage services (Section 10).

## 2. The 21st century vision of computing

With the advancement of modern society, basic essential services (*utilities*) are commonly provided such that everyone can easily obtain access to them. Today, utility services, such as water, electricity, gas, and telephony are deemed necessary for fulfilling daily life routines. These utility services are accessed so frequently that they need to be available whenever the consumer requires them at any time. Consumers are then able to pay service providers based on their usage of these utility services.

In 1969, Leonard Kleinrock [1], one of the chief scientists of the original Advanced Research Projects Agency Network (ARPANET) project which seeded the Internet, said: “As of now, computer networks are still in their infancy, but as they grow up and become sophisticated, we will probably see the spread of ‘*computer utilities*’ which, like present electric and telephone utilities, will service individual homes and offices across the country”. This vision of the computing utility based on the service provisioning model anticipates the massive transformation of the entire computing industry in the 21st century whereby computing services will be readily available on demand, like other utility services available in today’s society. Similarly, computing service users (consumers) need to pay providers only when they access computing services. In addition, consumers no longer need to invest heavily or encounter difficulties in building and maintaining complex IT infrastructure. Hence, software practitioners are facing numerous new challenges toward creating software for millions of consumers to use as a service, rather than to run on their individual computers.

The creation of the Internet has marked the foremost milestone towards achieving this grand 21st century vision of ‘*computer utilities*’ by forming a worldwide system of computer networks that enables individual computers to communicate with any other computers located elsewhere in the world. This internetworking of standalone computers reveals the promising potential of utilizing seemingly endless amount of distributed computing resources owned by various owners. As such, over the recent years, new computing paradigms (shown in Fig. 1) have been proposed and adopted to edge closer toward achieving this grand vision. Applications making use of these utility-oriented computing systems emerge simply as catalysts or market makers, which brings buyers and sellers together. This creates several trillion dollars worth of the utility/pervasive computing industry as noted by Sun Microsystems co-founder Bill Joy [2]. He also indicated “It would take time until these markets to mature to generate this kind of value. Predicting now which companies will capture the value is impossible. Many of them have not even been created yet.”

Grid computing [3] enables the sharing, selection, and aggregation of a wide variety of geographically distributed resources including supercomputers, storage systems, data sources, and specialized devices owned by different organizations for solving large-scale resource-intensive problems in science, engineering, and commerce. Inspired by the electrical power Grid’s pervasiveness, ease of use, and reliability [4], the motivation of Grid computing was initially driven by large-scale, resource (computational and data)-intensive scientific applications that required more resources than a single computer (PC, workstation, supercomputer) could have provided in a single administrative domain. Due to its potential to make impact on the 21st century as much as the electric power Grid did on the 20th century, Grid computing has been hailed as the next revolution after the Internet and the World Wide Web.

Peer-to-Peer (P2P) computing allows peer nodes (computers) to share content directly with one another in a decentralized manner. In pure P2P computing, there is no notion of clients or servers since all peer nodes are equal and concurrently be both clients and

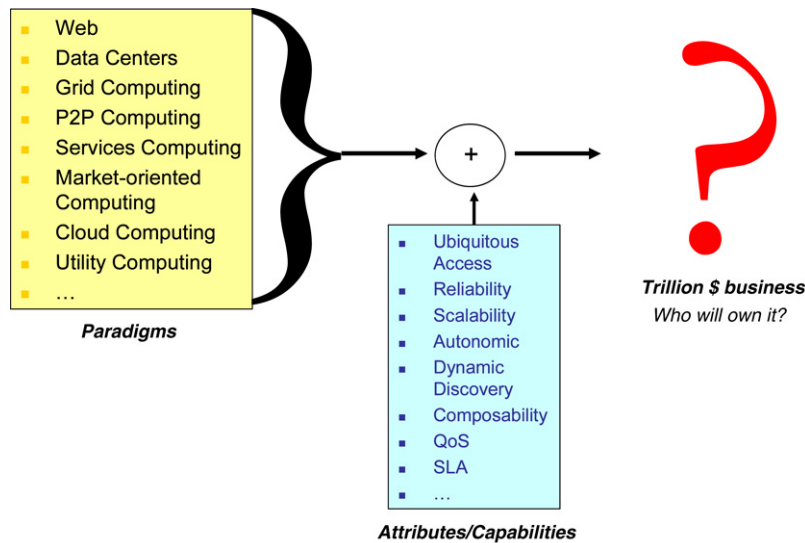


Fig. 1. Various paradigms promising to deliver IT as services.

servers. The goals of P2P computing include cost sharing or reduction, resource aggregation and interoperability, improved scalability and reliability, increased autonomy, anonymity or privacy, dynamism, and ad-hoc communication and collaboration [5].

Services computing focuses on the linkage between business processes and IT services so that business processes can be seamlessly automated using IT services. Examples of services computing technologies include Service-Oriented Architecture (SOA) and Web Services. The SOA facilitates interoperable services between distributed systems to communicate and exchange data with one another, thus providing a uniform means for service users and providers to discover and offer services respectively. The Web Services provides the capability for self-contained business functions to operate over the Internet.

Market-oriented computing views computing resources in economic terms such that resource users will need to pay resource providers for utilizing the computing resources [6]. Therefore, it is able to provide benefits, such as offering incentive for resource providers to contribute their resources for others to use and profit from it, regulating the supply and demand of computing resources at market equilibrium, offering incentive for resource users to back off when necessary, removing the need for a central coordinator (during the negotiation between the user and provider for establishing quality of service expectations and service pricing), and enabling both users and providers to make independent decisions to maximize their utility and profit respectively.

Today, the latest paradigm to emerge is that of Cloud computing [7] which promises reliable services delivered through next-generation data centers that are built on virtualized compute and storage technologies. Consumers will be able to access applications and data from a “Cloud” anywhere in the world on demand. The consumers are assured that the Cloud infrastructure is very robust and will always be available at any time. Computing services need to be highly reliable, scalable, and autonomic to support ubiquitous access, dynamic discovery and composability. In particular, consumers indicate the required service level through Quality of Service (QoS) parameters, which are noted in SLAs established with providers. Of all these paradigms, the recently emerged Cloud computing paradigm appears to be the most promising one to leverage and build on the developments from other paradigms.

### 3. Definitions, characteristics, and trends

In order to facilitate a clear understanding of what exactly is Cloud computing, we compare Cloud computing with two other recent, widely-adopted or explored computing paradigms: Cluster Computing and Grid Computing. We first examine the respective definitions of these three paradigms, then differentiate their specific characteristics, and finally highlight their recent web search trends.

#### 3.1. Definitions

A number of computing researchers and practitioners have attempted to define clusters, Grids, and Clouds [8] in various ways. Here are some definitions that we think are generic enough to stand the test of time.

The essence of Pfister’s [9] and Buyya’s [10] work defines clusters as follows:

- “A cluster is a type of parallel and distributed system, which consists of a collection of inter-connected stand-alone computers working together as a single integrated computing resource.”

Buyya defined one of the popular definitions for Grids at the 2002 Grid Planet conference, San Jose, USA as follows:

- “A Grid is a type of parallel and distributed system that enables the sharing, selection, and aggregation of geographically distributed ‘autonomous’ resources dynamically at runtime depending on their availability, capability, performance, cost, and users’ quality-of-service requirements.”

Based on our observation of the essence of what Clouds are promising to be, we propose the following definition:

- “A Cloud is a type of parallel and distributed system consisting of a collection of inter-connected and virtualized computers that are dynamically provisioned and presented as one or more unified computing resource(s) based on service-level agreements established through negotiation between the service provider and consumers.”

At a cursory glance, Clouds appear to be a combination of clusters and Grids. However, this is not the case. Clouds are clearly

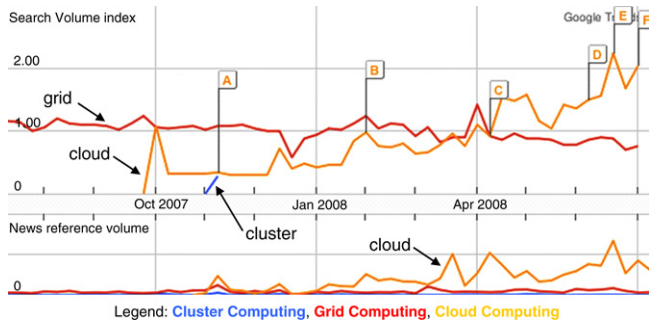


Fig. 2. Google search trends for the last 12 months.

next-generation data centers with nodes “virtualized” through hypervisor technologies such as VMs, dynamically “provisioned” on demand as a personalized resource collection to meet a specific service-level agreement, which is established through a “negotiation” and accessible as a composable service via Web Service technologies such as SOAP and REST.

### 3.2. Characteristics

A set of characteristics that helps distinguish cluster, Grid and Cloud computing systems is listed in Table 1. The resources in clusters are located in a single administrative domain and managed by a single entity whereas, in Grid systems, resources are geographically distributed across multiple administrative domains with their own management policies and goals. Another key difference between cluster and Grid systems arises from the way application scheduling is performed. The *schedulers* in cluster systems focus on enhancing the overall system performance and utility as they are responsible for the whole system. On the other hand, the schedulers in Grid systems called *resource brokers*, focusing on enhancing the performance of a specific application in such a way that its end-users’ QoS requirements are met.

Cloud computing platforms possess characteristics of both clusters and Grids, with its own special attributes and capabilities such strong support for virtualization, dynamically composable services with Web Service interfaces, and strong support for creating 3rd party, value added services by building on Cloud compute, storage, and application services. Thus, Clouds are promising to provide services to users without reference to the infrastructure on which these are hosted.

### 3.3. Web search trends

The popularity of different paradigms varies with time. The web search popularity, as measured by the Google search trends during the last 12 months, for terms “cluster computing”, “Grid computing”, and “Cloud computing” is shown in Fig. 2. From the Google trends, it can be observed that cluster computing was a popular term during 1990s, from early 2000 Grid computing become popular, and recently Cloud computing started gaining popularity.

Spot points in Fig. 2 indicate the release of news related to Cloud computing as follows:

- A IBM Introduces ‘Blue Cloud’ Computing, CIO Today – Nov 15 2007.
- B IBM, EU Launch RESERVOIR Research Initiative for Cloud Computing, IT News Online – Feb 7 2008.
- C Google and Salesforce.com in Cloud computing deal, Siliconpublic.com – Apr 14 2008.
- D Demystifying Cloud Computing, Intelligent Enterprise – Jun 11 2008.
- E Yahoo realigns to support Cloud computing, ‘core strategies’, San Antonio Business Journal – Jun 27 2008.

F Merrill Lynch Estimates “Cloud Computing” To Be \$100 Billion Market, SYS-CON Media – Jul 8 2008.

Other more recent news includes the following:

- Yahoo, Intel and HP form Cloud computing labs, Reseller News – Jul 29 2008.
- How Cloud Computing Is Changing The World, Pittsburgh Channel.com – Aug 4 2008.
- SIMtone Corporation Takes Cloud Computing to the Next Level with Launch of First Wireless, “Zero-Touch” Universal Cloud Computing Terminal, TMCnet – Sep 8 2008.

## 4. Market-oriented Cloud architecture

As consumers rely on Cloud providers to supply more of their computing needs, they will require specific QoS to be maintained by their providers in order to meet their objectives and sustain their operations. Cloud providers will need to consider and meet different QoS parameters of each individual consumer as negotiated in specific SLAs. To achieve this, Cloud providers can no longer continue to deploy traditional system-centric resource management architecture that do not provide incentives for them to share their resources and still regard all service requests to be of equal importance. Instead, market-oriented resource management [11,12] is necessary to regulate the supply and demand of Cloud resources to achieve market equilibrium (where supply = demand), providing feedback in terms of economic incentives for both Cloud consumers and providers, and promoting QoS-based resource allocation mechanisms that differentiate service requests based on their utility. In addition, clients can benefit from the “potential” cost reduction of providers, which could lead to a more competitive market and thus lower prices.

Fig. 3 shows the high-level architecture for supporting market-oriented resource allocation in Data Centers and Clouds. There are basically four main entities involved:

- **Users/Brokers:** Users or brokers acting on their behalf submit service requests from anywhere in the world to the Data Center and Cloud to be processed.
- **SLA Resource Allocator:** The SLA Resource Allocator acts as the interface between the Data Center/Cloud service provider and external users/brokers. It requires the interaction of the following mechanisms to support SLA-oriented resource management:
  - Service Request Examiner and Admission Control: When a service request is first submitted, the Service Request Examiner and Admission Control mechanism interprets the submitted request for QoS requirements before determining whether to accept or reject the request. Thus, it ensures that there is no overloading of resources whereby many service requests cannot be fulfilled successfully due to limited resources available. It also needs the latest status information regarding resource availability (from the VM Monitor mechanism) and workload processing (from the Service Request Monitor mechanism) in order to make resource allocation decisions effectively. Then, it assigns requests to VMs and determines resource entitlements for allocated VMs.
  - Pricing: The Pricing mechanism decides how service requests are charged. For instance, requests can be charged based on submission time (peak/off-peak), pricing rates (fixed/changing) or availability of resources (supply/demand). Pricing serves as a basis for managing the supply and demand of computing resources within the Data Center and facilitates in prioritizing resource allocations effectively.
  - Accounting: The Accounting mechanism maintains the actual usage of resources by requests so that the final cost can

**Table 1**  
Key characteristics of clusters, Grids, and Cloud systems.

Characteristics	Systems		
	Clusters	Grids	Clouds
Population	Commodity computers	High-end computers (servers, clusters)	Commodity computers and high-end servers and network attached storage
Size/scalability	100s	1000s	100s to 1000s
Node Operating System (OS)	One of the standard OSs (Linux, Windows)	Any standard OS (dominated by Unix)	A hypervisor (VM) on which multiple OSs run
Ownership	Single	Multiple	Single
Interconnection network/speed	Dedicated, high-end with low latency and high bandwidth	Mostly Internet with high latency and low bandwidth	Dedicated, high-end with low latency and high bandwidth
Security/privacy	Traditional login/password-based. Medium level of privacy – depends on user privileges.	Public/private key pair based authentication and mapping a user to an account. Limited support for privacy.	Each user/application is provided with a virtual machine. High security/privacy is guaranteed. Support for setting per-file access control list (ACL).
Discovery	Membership services	Centralised indexing and decentralised info services	Membership services
Service negotiation	Limited	Yes, SLA based	Yes, SLA based
User management	Centralised	Decentralised and also virtual organization (VO)-based	Centralised or can be delegated to third party
Resource management	Centralized	Distributed	Centralized/Distributed
Allocation/scheduling	Centralised	Decentralised	Both centralised/decentralised
Standards/inter-operability	Virtual Interface Architecture (VIA)-based	Some Open Grid Forum standards	Web Services (SOAP and REST)
Single system image	Yes	No	Yes, but optional
Capacity	Stable and guaranteed	Varies, but high	Provisioned on demand
Failure management (Self-healing)	Limited (often failed tasks/applications are restarted).	Limited (often failed tasks/applications are restarted).	Strong support for failover and content replication. VMs can be easily migrated from one node to other.
Pricing of services	Limited, not open market	Dominated by public good or privately assigned	Utility pricing, discounted for larger customers
Internetworking	Multi-clustering within an Organization	Limited adoption, but being explored through research efforts such as Gridbus InterGrid	High potential, third party solution providers can loosely tie together services of different Clouds
Application drivers	Science, business, enterprise computing, data centers	Collaborative scientific and high throughput computing applications	Dynamically provisioned legacy and web applications, Content delivery
Potential for building 3rd party or value-added solutions	Limited due to rigid architecture	Limited due to strong orientation for scientific computing	High potential – can create new services by dynamically provisioning of compute, storage, and application services and offer as their own isolated or composite Cloud services to users

be computed and charged to the users. In addition, the maintained historical usage information can be utilized by the Service Request Examiner and Admission Control mechanism to improve resource allocation decisions.

- o VM Monitor: The VM Monitor mechanism keeps track of the availability of VMs and their resource entitlements.
- o Dispatcher: The Dispatcher mechanism starts the execution of accepted service requests on allocated VMs.
- o Service Request Monitor: The Service Request Monitor mechanism keeps track of the execution progress of service requests.
- **VMs:** Multiple VMs can be started and stopped on-demand on a single physical machine to meet accepted service requests, hence providing maximum flexibility to configure various partitions of resources on the same physical machine to different specific requirements of service requests. In addition, multiple VMs can concurrently run applications based on different operating system environments on a single physical machine since every VM is completely isolated from one another on the same physical machine.
- **Physical Machines:** The Data Center comprises multiple computing servers that provide resources to meet service demands.

In the case of a Cloud as a commercial offering to enable crucial business operations of companies, there are critical QoS parameters to consider in a service request, such as time, cost, reliability and trust/security. In particular, QoS requirements cannot be static and may change over time due to continuing changes in business operations and operating environments. In short, there should be greater importance on customers since they

pay for accessing services in Clouds. In addition, the state-of-the-art in Cloud computing has no or limited support for dynamic negotiation of SLAs between participants and mechanisms for automatic allocation of resources to multiple competing requests. Recently, we have developed negotiation mechanisms based on alternate offers protocol for establishing SLAs [13]. These have high potential for their adoption in Cloud computing systems built using VMs.

Commercial offerings of market-oriented Clouds must be able to:

- Support customer-driven service management based on customer profiles and requested service requirements,
- Define computational risk management tactics to identify, assess, and manage risks involved in the execution of applications with regards to service requirements and customer needs,
- Derive appropriate market-based resource management strategies that encompass both customer-driven service management and computational risk management to sustain SLA-oriented resource allocation,
- Incorporate autonomic resource management models that effectively self-manage changes in service requirements to satisfy both new service demands and existing service obligations, and
- Leverage VM technology to dynamically assign resource shares according to service requirements.

## 5. Resource management strategies for market-oriented Clouds

Since customer satisfaction is the crucial success factor to excel in the service industry [14], computing service providers

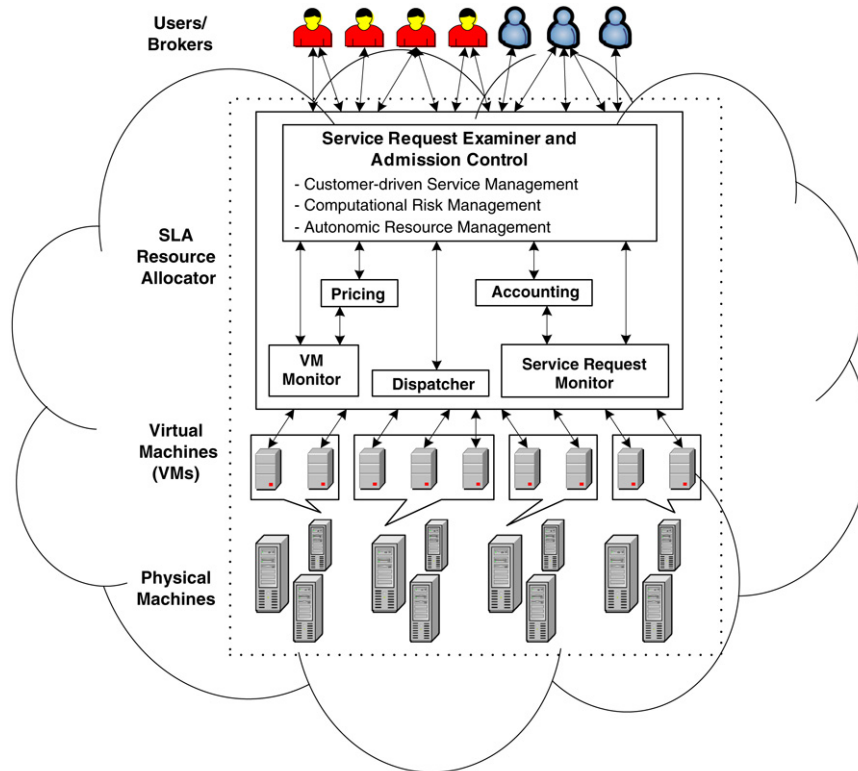


Fig. 3. High-level market-oriented Cloud architecture.

have to be aware of user-centric objectives and meet them in order to achieve customer satisfaction. But, many service quality factors can influence customer satisfaction [14,15]. Hence, we need to design SLA-oriented resource management strategies for Data Centers and Clouds that provide personalized attention to customers, such as enabling communication to keep customers informed and obtain feedback from them, increasing access and approachability to customers, and understanding specific needs of customers. These strategies can also encourage trust and confidence in customers by emphasizing on the security measures undertaken against risks and doubts, the credibility of the provider, and the courtesy towards customers.

Our initial work [16] has also presented examples of how various elements of utility-based resource management can be perceived as risks and hence identified risk analysis from the field of economics as a probable solution to evaluate them. However, the entire risk management process [17,18] comprises of many steps and needs to be studied thoroughly so as to fully realize its effectiveness in managing risks. Hence, we need to first establish the context of risk management in Data Centers and Clouds, and then identify the risks involved. Each of the identified risks will be thoroughly assessed, before deriving appropriate strategies to manage these risks.

In addition, service requirements of users can change over time and thus may require amendments of original service requests. As such, our proposed resource management strategies will be able to self-manage the reservation process continuously by monitoring current service requests, amending future service requests, and adjusting schedules and prices for new and amended service requests accordingly. Hence, we need to investigate self-configuring components to satisfy new service requirements, so that more autonomic and intelligent Data Centers and Clouds can better manage the limited supply of resources with dynamically changing service demand. For users, there can be brokering systems acting on their behalf to select suitable providers

and negotiate with them to achieve ideal service contracts. Thus, providers also require autonomic resource management to selectively choose appropriate requests to accept and execute depending on a number of operating factors, such as the expected availability and demand of services (both current and future), and existing service obligations.

Recently, virtualization [19] has enabled the abstraction of computing resources such that a single physical machine is able to function as a set of multiple logical VMs. A key benefit of VMs is the ability to host multiple operating system environments which are completely isolated from one another on the same physical machine. Another benefit is the capability to configure VMs to utilize different partitions of resources on the same physical machine. For example, on a physical machine, one VM can be allocated 10% of the processing power, while another VM can be allocated 20% of the processing power. Hence, we need to leverage existing VM technologies so that VMs can be started and stopped dynamically to meet the changing demand of resources by users as opposed to limited resources on a physical machine. In particular, we need to investigate how VMs can be assigned various resource management policies catering to different user needs and demands to better support the implementation of SLA-oriented resource allocation for Data Centers and Clouds.

## 6. Global cloud exchanges and markets

Enterprises currently employ Cloud services in order to improve the scalability of their services and to deal with bursts in resource demands. However, at present, service providers have inflexible pricing, generally limited to flat rates or tariffs based on usage thresholds, and consumers are restricted to offerings from a single provider at a time. Also, many providers have proprietary interfaces to their services thus restricting the ability of consumers to swap one provider for another.

For Cloud computing to mature, it is required that the services follow standard interfaces. This would enable services to be

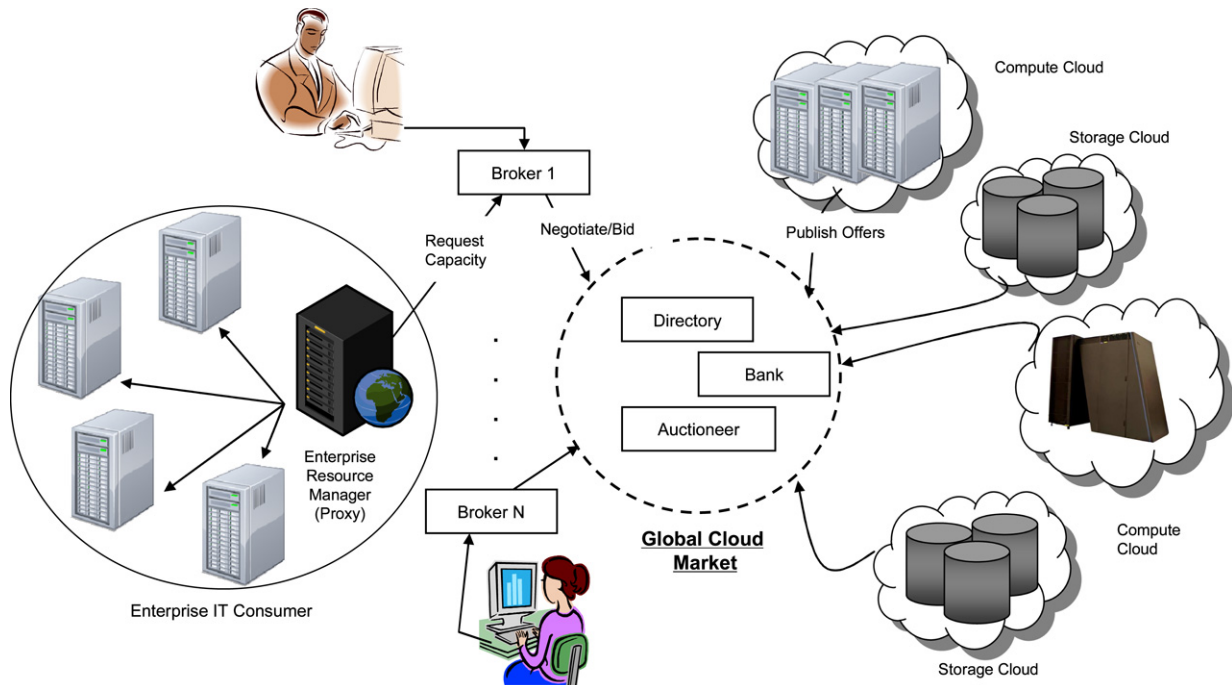


Fig. 4. Global Cloud exchange and market infrastructure for trading services.

commoditized and thus, would pave the way for the creation of a market infrastructure for trading in services. An example of such a market system, modeled on real-world exchanges, is shown in Fig. 4. The market directory allows participants to locate providers or consumers with the right offers. Auctioneers periodically clear bids and asks received from market participants. The banking system ensures that financial transactions pertaining to agreements between participants are carried out.

Brokers perform the same function in such a market as they do in real-world markets: they mediate between consumers and providers by buying capacity from the provider and sub-leasing these to the consumers. A broker can accept requests from many users who have a choice of submitting their requirements to different brokers. Consumers, brokers and providers are bound to their requirements and related compensations through SLAs. An SLA specifies the details of the service to be provided in terms of metrics agreed upon by all parties, and penalties for meeting and violating the expectations, respectively.

Such markets can bridge disparate Clouds allowing consumers to choose a provider that suits their requirements by either executing SLAs in advance or by buying capacity on the spot. Providers can use the markets in order to perform effective capacity planning. A provider is equipped with a price-setting mechanism which sets the current price for the resource based on market conditions, user demand, and current level of utilization of the resource. Pricing can be either fixed or variable depending on the market conditions. An admission-control mechanism at a provider's end selects the auctions to participate in or the brokers to negotiate with, based on an initial estimate of the utility. The negotiation process proceeds until an SLA is formed or the participants decide to break off. These mechanisms interface with the resource management systems of the provider in order to guarantee the allocation being offered or negotiated can be reclaimed, so that SLA violations do not occur. The resource management system also provides functionalities such as advance reservations that enable guaranteed provisioning of resource capacity.

Brokers gain their utility through the difference between the price paid by the consumers for gaining resource shares and that

paid to the providers for leasing their resources. Therefore, a broker has to choose those users whose applications can provide it maximum utility. A broker interacts with resource providers and other brokers to gain or to trade resource shares. A broker is equipped with a negotiation module that is informed by the current conditions of the resources and the current demand to make its decisions.

Consumers have their own utility functions that cover factors such as deadlines, fidelity of results, and turnaround time of applications. They are also constrained by the amount of resources that they can request at any time, usually by a limited budget. Consumers also have their own limited IT infrastructure that is generally not completely exposed to the Internet. Therefore, a consumer participates in the utility market through a resource management proxy that selects a set of brokers based on their offerings. He then forms SLAs with the brokers that bind the latter to provide the guaranteed resources. The enterprise consumer then deploys his own environment on the leased resources or uses the provider's interfaces in order to scale his applications.

The idea of utility markets for computing resources has been around for a long time. Recently, many research projects such as SHARP [20], Tycoon [21], Bellagio [22], and Shirako [23] have come up with market structures for trading in resource allocations. These have particularly focused on trading in VM-based resource slices on networked infrastructures such as PlanetLab. The Gridbus project has created a resource broker that is able to negotiate with resource providers. Thus, the technology for enabling utility markets is already present and ready to be deployed.

However, significant challenges persist in the universal application of such markets. Enterprises currently employ conservative IT strategies and are unwilling to shift from the traditional controlled environments. Cloud computing uptake has only recently begun and many systems are in the proof-of-concept stage. Regulatory pressures also mean that enterprises have to be careful about where their data gets processed, and therefore, are not able to employ Cloud services from an open market. This could be mitigated through SLAs that specify strict constraints on the location of the resources. However, another open issue is how the participants in such a market can obtain restitution in case an SLA is violated. This

**Table 2**  
Comparison of some representative Cloud platforms.

Property	System				
	Amazon Elastic compute cloud (EC2)	Google App engine	Microsoft Azure	Sun Network.com (Sun Grid)	GRIDS Lab Aneka
Focus	Infrastructure	Platform	Platform	Infrastructure	Software platform for enterprise Clouds
Service type	Compute, storage (Amazon S3)	Web application	Web and non-web application	Compute	Compute
Virtualization	OS level running on a Xen hypervisor	Application container	OS level through fabric controller	Job management system (Sun Grid Engine)	Resource manager and scheduler
Dynamic negotiation of QoS parameters	None	None	None	None	SLA-based resource reservation on Aneka side.
User access interface	Amazon EC2 command-line tools	Web-based administration console	Microsoft windows azure portal	Job submission scripts, Sun Grid web portal	Workbench, web-based portal
Web APIs	Yes	Yes	Yes	Yes	Yes
Value-added service providers	Yes	No	Yes	Yes	No
Programming framework	Customizable linux-based Amazon Machine Image (AMI)	Python	Microsoft .NET	Solaris OS, Java, C, C++, FORTRAN	APIs supporting different programming models in C# and other .Net supported languages

motivates the need for a legal framework for agreements in such markets, a research issue that is out of scope of themes pursued in this paper.

## 7. Emerging Cloud platforms

Industry analysts have made bullish projections on how Cloud computing will transform the entire computing industry. According to a recent Merrill Lynch research note [24], Cloud computing is expected to be a “\$160-billion addressable market opportunity, including \$95-billion in business and productivity applications, and another \$65-billion in online advertising”. Another research study by Morgan Stanley [25] has also identified Cloud computing as one of the prominent technology trends. As the computing industry shifts toward providing Platform as a Service (PaaS) and Software as a Service (SaaS) for consumers and enterprises to access on demand regardless of time and location, there will be an increase in the number of Cloud platforms available. Recently, several academic and industrial organizations have started investigating and developing technologies and infrastructure for Cloud Computing. Academic efforts include Virtual Workspaces [26], OpenNebula [27], and Reservoir [28]. In this section, we compare six representative Cloud platforms with industrial linkages in Table 2.

Amazon Elastic Compute Cloud (EC2) [29] provides a virtual computing environment that enables a user to run Linux-based applications. The user can either create a new Amazon Machine Image (AMI) containing the applications, libraries, data and associated configuration settings, or select from a library of globally available AMIs. The user then needs to upload the created or selected AMIs to Amazon Simple Storage Service (S3), before he can start, stop, and monitor instances of the uploaded AMIs. Amazon EC2 charges the user for the time when the instance is alive, while Amazon S3 [30] charges for any data transfer (both upload and download).

Google App Engine [31] allows a user to run web applications written using the Python programming language. Other than supporting the Python standard library, Google App Engine also supports Application Programming Interfaces (APIs) for the datastore, Google Accounts, URL fetch, image manipulation, and email services. Google App Engine also provides a web-based Administration Console for the user to easily manage his running web applications. Currently, Google App Engine is free to use with up to 500MB of storage and about 5 million page views per month.

Microsoft Azure [32] aims to provide an integrated development, hosting, and control Cloud computing environment so that software developers can easily create, host, manage, and scale both Web and non-web applications through Microsoft data centers. To achieve this aim, Microsoft Azure supports a comprehensive collection of proprietary development tools and protocols which consists of Live Services, Microsoft .NET Services, Microsoft SQL Services, Microsoft SharePoint Services, and Microsoft Dynamics CRM Services. Microsoft Azure also supports Web APIs such as SOAP and REST to allow software developers to interface between Microsoft or non-Microsoft tools and technologies.

Sun network.com (Sun Grid) [33] enables the user to run Solaris OS, Java, C, C++, and FORTRAN based applications. First, the user has to build and debug his applications and runtime scripts in a local development environment that is configured to be similar to that on the Sun Grid. Then, he needs to create a bundled zip archive (containing all the related scripts, libraries, executable binaries and input data) and upload it to Sun Grid. Finally, he can execute and monitor the application using the Sun Grid web portal or API. After the completion of the application, the user will need to download the execution results to his local development environment for viewing.

Aneka [34], which is being commercialized through Manjrasoft, is a .NET-based service-oriented resource management platform. It is designed to support multiple application models, persistence and security solutions, and communication protocols such that the preferred selection can be changed at anytime without affecting an existing Aneka ecosystem. To create an Aneka Cloud, the service provider only needs to start an instance of the configurable Aneka container hosting required services on each selected desktop computer. The purpose of the Aneka container is to initialize services and acts as a single point for interaction with the rest of the Aneka Cloud. Aneka provides SLA support such that the user can specify QoS requirements such as deadline (maximum time period which the application needs to be completed in) and budget (maximum cost that the user is willing to pay for meeting the deadline). The user can access the Aneka Cloud remotely through the Gridbus broker. The Gridbus broker [35] also enables the user to negotiate and agree upon the QoS requirements to be provided by the service provider.

## 8. Aneka: From enterprise Grids to market-oriented Cloud computing

We are working towards implementing a market-oriented Cloud using a .NET-based service-oriented resource management

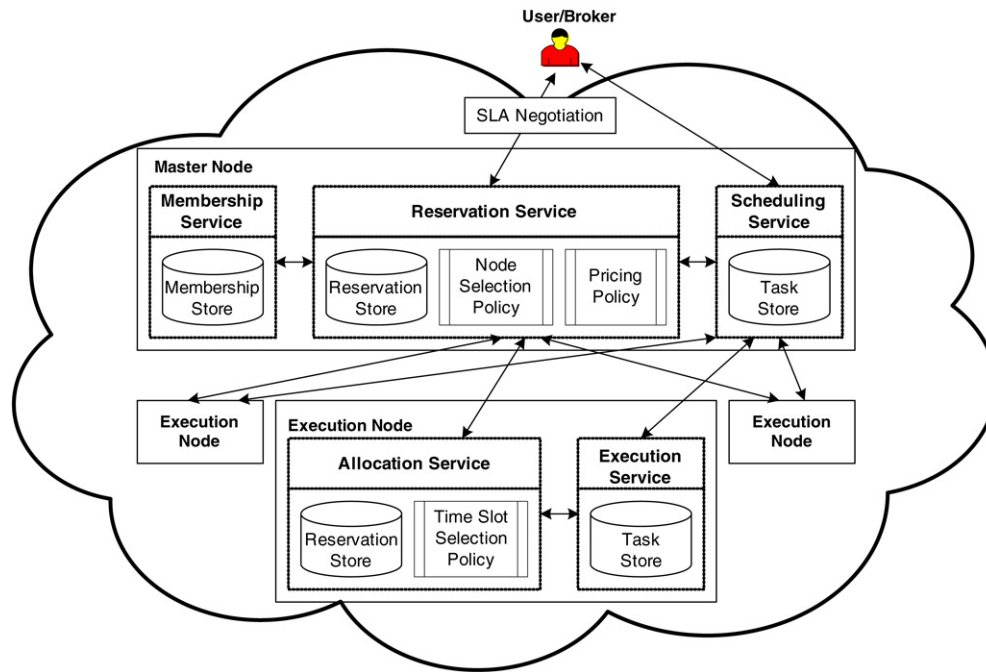


Fig. 5. Interaction of services in Aneka Cloud environment.

platform called Aneka [34]. Aneka is initially developed as a 3rd generation enterprise Grid technology. Recently, various new capabilities have been added to exhibit properties and potentials of the Cloud computing paradigm. An enterprise Grid [36] harnesses computing resources of desktop computers (connected over an internal network or the Internet) within an enterprise without affecting the productivity of their users. Hence, it increases the amount of computing resources available within an enterprise to accelerate application performance. This capability can be combined with other dedicated resources in the enterprise to enhance the overall system capability and reliability.

To support scalability, the Aneka container is designed to be lightweight by providing the bare minimum functionality needed for an Aneka Cloud node. It provides the base infrastructure that consists of services for persistence, security (authorization, authentication and auditing), and communication (message handling and dispatching). The Aneka container can host any number of optional services that can be added to augment the capabilities of an Aneka Cloud node. Examples of optional services are indexing, scheduling, execution, and storage services. This provides a single, flexible and extensible framework for orchestrating various application models. This section describes how pricing can be implemented in an Aneka Cloud with advanced reservations.

### 8.1. Market-oriented resource pricing and allocation in Aneka Cloud

To create an Aneka Cloud, we implement a bi-hierarchical advance reservation mechanism with a Reservation Service at a master node that coordinates multiple execution nodes and an Allocation Service at each execution node that keeps track of the reservations at that node. This architecture was previously introduced in [13]. To use the Aneka Cloud, the resource user (or a broker acting on its behalf) first makes advanced reservations during the reservation phase for resources required at a designated time in the future. If the reservation phase is successful, the user/broker can then submit applications later during the execution phase when the designated time in the future arrives.

Fig. 5 shows that the process of allocating advanced reservations happens in two levels: the Allocation Service at each execution node and the Reservation Service at the master node. Both

services are designed to support pluggable policies so that the provider has the flexibility to easily customize and replace existing policies for different levels and/or nodes without interfering with the overall resource management architecture.

During the reservation phase, the user/broker submits reservation requests through the Reservation Service at the master node. The Reservation Service discovers available execution nodes in the Aneka Cloud by interacting with the Allocation Service on them. The Allocation Service at each execution node keeps track of all reservations that have been confirmed for the node and can thus check whether a new request can be satisfied or not.

The Allocation Service determines how to schedule a new reservation at the execution node. For simplicity, we implement the same time slot selection policy for the Allocation Service at every execution node. The Allocation Service allocates the requested time slot if the slot is available. Otherwise, it assigns the next available time slot after the requested start time that can meet the required duration. By allocating reservations at each execution node instead of at the master node, computation overheads arising from making allocation decisions are distributed across multiple nodes and thus minimized, as compared to overhead accumulation at a single master node.

The Reservation Service performs node selection by choosing the required number of available time slots from execution nodes and administers admission control by accepting or rejecting a reservation request. It also calculates the price for a confirmed reservation based on the implemented pricing policy. Available time slots are selected taking into account the application requirement of the user.

The application requirement considered here is the task parallelism to execute an application. A sequential application has a single task and thus needs a single processor to run, while a parallel application needs a required number of processors to concurrently run at the same time.

For a sequential application, the selected time slots need not have the same start and end times. Hence, available time slots with the lowest prices are selected first. If there are multiple available time slots with the same price, then those with the earliest start time available are selected first. This ensures that

the cheapest requested time slot is allocated first if it is available. Selecting available time slots with the lowest prices first is fair and realistic. In reality, reservations that are confirmed earlier enjoy the privilege of cheaper prices, as compared to reservation requests that arrive later.

But, for a parallel application, all the selected time slots must have the same start and end times. Again, the earliest time slots (with the same start and end times) are allocated first to ensure the requested time slot is allocated first if available. If there are more available time slots (with the same start and end times) than the required number of time slots, then those with the lowest prices are selected first.

The admission control operates according to the service requirement of the user. The service requirements examined are the deadline and budget to complete an application. We assume both deadline and budget are hard constraints. Hence, a confirmed reservation must not end after the deadline and cost more than the budget. Therefore, a reservation request is not accepted if there is insufficient number of available time slots on execution nodes that ends within the deadline and the total price of the reservation costs more than the budget.

During the execution phase, the user/broker submits applications to be executed to the Scheduling Service at the master node. The Scheduling Service first checks whether the submission satisfies the starting time, ending time, and duration of the reservation that have been specified by the user/broker during the reservation phase. If the reservation is still valid, the Scheduling Service then determines whether any of the reserved execution nodes are available before dispatching applications to them for execution, otherwise applications are queued to wait for the next available execution nodes that are part of the reservation. The Execution Service at each execution node starts executing an application after receiving it from the Scheduling Service and updates the Scheduling Service of changes in execution status. Hence, the Scheduling Service can monitor executions for an application and notify the user/broker upon completion.

## 8.2. Performance evaluation

High-end computing systems such as Clouds are used for hosting applications containing short-lived or long-lived processing tasks. Applications providing Internet services often have short-lived tasks and are designed to serve millions of users concurrently. Examples include search engines (e.g. Google), e-commerce sites (e.g. Amazon.com online shopping store), and social networking sites (e.g. Facebook). Many business and scientific applications such as investment risk analysis, supply chain management, flight simulation, and molecular docking often contain tasks that are resource-intensive and long-lived. We will first present the performance of HPC workload (with long-lived tasks) in our Aneka enterprise Cloud environment. We will then discuss another performance study on Internet-based services workload (with short-lived tasks).

### 8.2.1. High Performance Computing (HPC) workload

Fig. 6 shows the Aneka enterprise Cloud environment setup used for performance evaluation. The Aneka Cloud contains 33 personal computers (PCs) with 1 master node and 32 execution nodes located across 3 student computer laboratories in the Department of Computer Science and Software Engineering, The University of Melbourne. This setup demonstrates that the Aneka Cloud is able to present a unified resource to the users/brokers by harnessing computing resources located physically apart in the 3 laboratories.

Synthetic workloads are created by utilizing trace data of HPC applications. The experiments utilize 238 reservation requests

**Table 3**  
Pricing mechanisms.

Name	Configured pricing parameters
FixedMax	\$3/CPU/Hr
FixedMin	\$1/CPU/Hr
FixedTimeMax	\$1/CPU/Hr (12AM–12PM) \$3/CPU/Hr (12PM–12AM)
FixedTimeMin	\$1/CPU/Hr (12AM–12PM) \$2/CPU/Hr (12PM–12AM)
Libra+\$Max	\$1/CPU/Hr ( $PBase_j$ ), $\alpha = 1$ , $\beta = 3$
Libra+\$Min	\$1/CPU/Hr ( $PBase_j$ ), $\alpha = 1$ , $\beta = 1$
Libra+\$Auto	same as Libra+\$Min

in the last 7 days of the SDSC SP2 trace (April 1998 to April 2000) version 2.2 from Feitelson's Parallel Workloads Archive [37]. The SDSC SP2 trace from the San Diego Supercomputer Center (SDSC) in USA is chosen due to the highest resource utilization of 83.2% among available traces to ideally model a heavy workload scenario. However, the trace only provides the inter-arrival times of reservation requests, the number of processors to be reserved (downscaled from a maximum of 128 nodes in the trace to a maximum of 32 nodes), and the duration to be reserved. Hence, we adopt a methodology similar to that adopted by Irwin et al. [38] to synthetically assign service requirements (deadline and budget) through two request classes: (i) low urgency and (ii) high urgency.

A reservation request  $i$  in the *low urgency* class has a deadline of high  $deadline_i/duration_i$  value and budget of low  $budget_i/f(duration_i)$  value.  $f(duration_i)$  is a function representing the minimum budget required based on  $duration_i$ . Conversely, each request  $i$  in the *high urgency* class has a deadline of low  $deadline_i/duration_i$  value and budget of high  $budget_i/f(duration_i)$  value. This is realistic since a user who submits a more urgent request to be met within a shorter deadline offers a higher budget for the short notice. Values are normally distributed within each of the deadline and budget parameters.

For simplicity, the Aneka Cloud only charges users for utilizing the computing resource type based on per processor (CPU) per hour (h). Thus, users are not charged for using other resource types such as memory, storage, and bandwidth. In addition, we assume that every user/broker can definitely accept another reservation time slot proposed by the Aneka Cloud if the requested one is not possible, provided that the proposed time slot still satisfies both application and service requirements of the user.

As listed in Table 3, we evaluate the performance of seven pricing mechanisms representing three basic types: (i) Fixed, (ii) FixedTime, and (iii) Libra+\$. Each of these three pricing mechanisms has a maximum and minimum configuration to highlight their performance range. The Fixed mechanism charges a fixed price at all times. The FixedTime mechanism charges a fixed price for different time periods of resource usage where a lower price is charged for off-peak (12AM–12PM) and a higher price for peak (12PM–12AM).

Libra+\$ [39] uses a more fine-grained pricing function that satisfies four essential requirements for pricing of resources to prevent workload overload: (i) flexible, (ii) fair, (iii) dynamic, and (iv) adaptive. The price  $P_{ij}$  for per unit of resource utilized by reservation request  $i$  at compute node  $j$  is computed as:  $P_{ij} = (\alpha * PBase_j) + (\beta * PUtil_{ij})$ . The base price  $PBase_j$  is a static pricing component for utilizing a resource at node  $j$  which can be used by the service provider to charge the minimum price so as to recover the operational cost. The utilization price  $PUtil_{ij}$  is a dynamic pricing component which is computed as a factor of  $PBase_j$  based on the utilization of the resource at node  $j$  for the required deadline of request  $i$ :  $PUtil_{ij} = RESMax_j/RESFree_{ij} * PBase_j$ .  $RESMax_j$  and  $RESFree_{ij}$  are the maximum units and remaining free units of the resource at node  $j$  for the deadline duration of request  $i$  respectively. Thus,  $RESFree_{ij}$  has been deducted units of resource

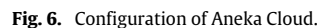


Fig. 6. Configuration of Aneka Cloud.

committed for other confirmed reservations and request  $i$  for its deadline duration.

The factors  $\alpha$  and  $\beta$  for the static and dynamic components of Libra+\$ respectively provides the flexibility for the service provider to easily configure and modify the weight of the static and dynamic components on the overall price  $P_{ij}$ . Libra+\$ is fair since requests are priced based on the amount of different resources utilized. It is also dynamic because the overall price of a request varies depending on the availability of resources for the required deadline. Finally, it is adaptive as the overall price is adjusted depending on the current supply and demand of resources to either encourage or discourage request submission.

However, these three mechanisms rely on static pricing parameters that are difficult to be accurately derived by the service provider to produce the best performance where necessary. Hence, we propose Libra+\$Auto [40], an autonomic Libra+\$ that automatically adjusts  $\beta$  based on the availability of compute nodes. Libra+\$Auto thus considers the pricing of resources across nodes, unlike Libra+\$ which only considers pricing of resources at each node  $j$  via  $P_{ij}$ .

Fig. 7 shows the normalized pricing and revenue performance of seven pricing mechanisms in the Aneka Cloud for high urgency requests (with short deadline and high budget) from sequential applications (requiring one processor to execute) over a 7-days time period. In Fig. 7, the two performance metrics are: (i) the price for a confirmed reservation (in \$/CPU/Hr) and (ii) the accumulated revenue for confirmed reservations (in \$). Both metrics have been normalized to produce standardized values within the range of 0 to 1 for easier comparison. The revenue of a confirmed reservation is the total sum of revenue across all its reserved nodes calculated using the assigned price (depending on the specific pricing mechanism) and the reserved duration at each node. Then, the price of a confirmed reservation is computed to reflect the average price across all its reserved nodes.

In Fig. 7, out of the four fixed pricing mechanisms listed in Table 3, FixedMax provides the highest revenue (maximum bound), followed by FixedTimeMax, FixedTimeMin, and FixedMin with the lowest revenue (minimum bound). Nevertheless, FixedTime mechanisms is easier to derive and more reliable than Fixed mechanisms since it supports a range of prices across various time periods of resource usage. But, all four mechanisms do not consider service requirements of users such as deadline and budget.

On the other hand, Libra+\$ charges a lower price for a request with longer deadline as an incentive to encourage users to submit requests with longer deadlines that are more likely to be accommodated than shorter deadlines. For a request with short deadline, Libra+\$Max and Libra+\$Min charge a higher price relative to their  $\beta$  in Table 3. Libra+\$Max provides higher revenue than Libra+\$Min due to a higher value of  $\beta$ .

Both Libra+\$Auto and Libra+\$Max are able to provide a significantly higher revenue than other pricing mechanisms through higher prices for shorter deadlines. Fig. 7 shows that Libra+\$Auto continues increasing prices to higher than that of Libra+\$Max and other pricing mechanisms when demand is high such as during the later half of day 1, 2, 3, and 5. But, when demand is low such as during the early half of day 2, 3, 5, and 6, Libra+\$Auto keeps reducing prices to lower than that of Libra+\$Max to accept requests that are not willing to pay more. Hence, Libra+\$Auto is able to exploit budget limits to achieve the highest revenue by automatically adjusting to a higher  $\beta$  to increase prices when the availability of nodes is low and a lower  $\beta$  to reduce prices when there are more unused nodes which will otherwise be wasted.

### 8.2.2. Internet-based services workload

MapReduce [41] is one of the most popular programming models designed for data centers. It was originally proposed by Google to handle large-scale web search applications and has been proved to be an effective programming model for













