**IBM**

# Seeding the Clouds: Key Infrastructure Elements for Cloud Computing

# Table of Contents

## Executive summary

Cloud computing is an emerging computing model by which users can gain access to their applications from anywhere, through any connected device. A user-centric interface makes the cloud infrastructure supporting the applications transparent to users. The applications reside in massively scalable data centers where computational resources can be dynamically provisioned and shared to achieve significant economies of scale. Thanks to a strong service management platform, the management costs of adding more IT resources to the cloud can be significantly lower than those associated with alternate infrastructures.

What is driving the adoption of cloud computing? Many factors, including the proliferation of smart mobile devices, high-speed connectivity, higher density computing and data-intensive Web 2.0 applications.

As a result, vendors across the IT industry have announced cloud computing efforts of varying capabilities and among corporate clients there is an increasing interest in aspects of the cloud, such as infrastructure outsourcing, software as a service key processes as a service and next-generation distributed computing.

IBM is uniquely positioned to help these clients adopt cloud computing technologies and management techniques to improve the efficiency and flexibility of their data centers. With proven expertise dating back to its pioneering position in the virtualization space during the 1960s, IBM has recently introduced its vision of a data center that supports a dynamic infrastructure. This vision brings together the strengths of the Web-centric cloud computing model and today's enterprise data center. It provides request-driven, dynamic allocation of computing resources for a mix of workloads on a massively scalable, heterogeneous and virtualized infrastructure. Furthermore, it is optimized for security, data integrity, resiliency and transaction processing. Thanks to extensive experience in both enterprise data centers and cloud computing, IBM is exceptionally well prepared to provide clients with the best solutions to achieve this vision.

IBM has been working with leading-edge clients around the world, such as Google and the Government of Wuxi in China, to define best practices for running data centers with workloads ranging from Web 2.0 applications to mission-critical transaction processing systems. Specifically, IBM has been defining and enhancing a cloud computing framework for running large scale data centers that enables key functionality for hosting a wide range of applications. This framework now includes automation for the complex, time-consuming processes of provisioning servers, networks, storage, operating systems and middleware. It also provides support for extremely data-intensive workloads and supports requirements for resiliency and security.

IBM's hands-on experience setting up cloud data centers represents a major step in the evolution of data centers in support of a dynamic infrastructure. This paper describes a high-level cloud computing infrastructure services framework and the underlying technology enablers, such as virtualization, automation, self-service portal, monitoring and capacity planning. It also discusses examples of, and the value propositions for, certain data centers that have been built in this manner to date. These data centers can host a mix of workloads, from Java™ 2 Enterprise Edition (J2EE) applications to software development to test environments to data-intensive business intelligence analytics.

## Introduction

### Business value of cloud computing

Cloud computing is both a business delivery model and an infrastructure management methodology. The business delivery model provides a user experience by which hardware, software and network resources are optimally leveraged to provide innovative services over the Web, and servers are provisioned in accordance with the logical needs of the service using advanced, automated tools. The cloud then enables the service creators, program administrators and others to use these services via a Web-based interface that abstracts away the complexity of the underlying dynamic infrastructure.

The infrastructure management methodology enables IT organizations to manage large numbers of highly virtualized resources as a single large resource. It also allows IT organizations to massively increase their data center resources without significantly increasing the number of people traditionally required to maintain that increase.

For organizations currently using traditional infrastructures, a cloud will enable users to consume IT resources in the data center in ways that were never available before. Companies that employ traditional data center management practices know that making IT resources available to an end user can be time-intensive. It involves many steps, such as procuring hardware; finding raised floor space and sufficient power and cooling; allocating administrators to install operating systems, middleware and software; provisioning the network; and securing the environment. Most companies find that this process can take upwards of two to three months. Those IT organizations that are re-provisioning existing hardware resources find that it still takes several weeks to accomplish. A cloud dramatically alleviates this problem by implementing automation, business workflows and resource abstraction that allows a user to browse a catalog of IT services, add them to a shopping cart and submit the order. After an administrator approves the order, the cloud does the rest. This process reduces the time required to make those resources available to the customer from months to minutes.

The cloud also provides a user interface that allows both the user and the IT administrator to easily manage the provisioned resources through the life cycle of the service request. After a user's resources have been delivered by a cloud, the user can track the order, which typically consists of some number of servers and software, and view the health of those resources; add servers; change the installed software; remove servers; increase or decrease the allocated processing power, memory or storage; and even start, stop and restart servers. These are self-service functions that can be performed 24 hours a day and take only minutes to perform. By contrast, in a non-cloud environment, it could take hours or days for someone to have a server restarted or hardware or software configurations changed.

The business model of a cloud facilitates more efficient use of existing resources. Clouds can require users to commit to predefined start and end dates for resource requests. This helps IT organizations to more efficiently repurpose resources that often get forgotten or go unused. When users realize they can get resources within minutes of a request, they are less likely to hoard resources that are otherwise very difficult to acquire.

Clouds provide request-driven, dynamic allocation of computing resources for a mix of workloads on a massively scalable, heterogeneous and virtualized infrastructure. The value of a fully automated provisioning process that is security compliant and automatically customized to user's needs results in:

- *Significantly reduced time to introduce technologies and innovations;*
- *Cost savings in labor for designing, procuring and building hardware and software platforms;*
- *Cost savings by avoiding human error in the configuration of security, networks and the software provisioning process; and*
- *Cost elimination through greater use and reuse of existing resources, resulting in better efficiency.*

The cloud computing model reduces the need for capacity planning at an application level. The user of an application can request resources from the cloud and obtain them in less than an hour. A user who needs more resources can submit another request and obtain more resources within minutes, and in a policy-based system, no interaction is needed at all; resource changes are performed dynamically. Thus, it is far less important to correctly predict the capacity requirements for an application than it is in traditional data centers, and capacity planning is simplified because it is performed only once for the entire data center.

Today's IT realities make cloud computing a good fit for meeting the needs of both *IT providers* (who demand unprecedented flexibility and efficiency, lower costs and complexity and support for varied and huge workloads) and *Internet users* (who expect availability, function and speed).

As technology such as virtualization and corresponding management services like automation, monitoring and capacity planning services become more mature, cloud computing will become more widely used for increasingly diverse and even mission-critical workloads.

### Evolution of cloud computing

*This section reviews the history of cloud computing and introduces the IBM vision for cloud computing that supports dynamic infrastructures. The following section introduces an infrastructure framework for a data center and discusses the virtualized environment and infrastructure management. Subsequently, existing cloud infrastructures and their applications are described.*

Cloud computing is an important topic. However, it is not a revolutionary new development, but an evolution that has taken place over several decades, as shown in Figure 1.
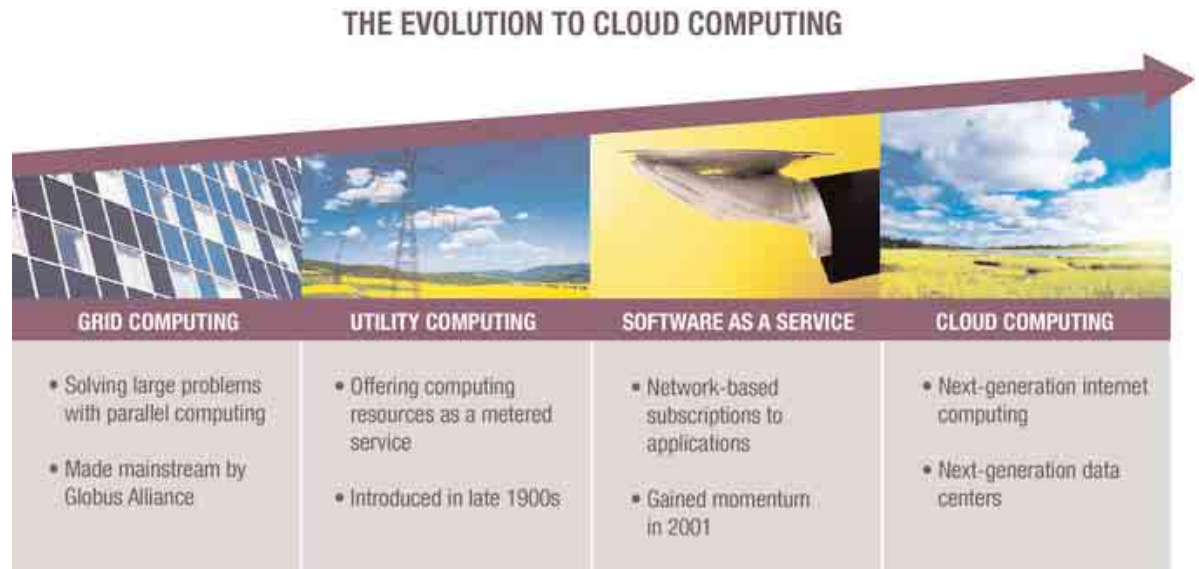


Figure 1. Evolution toward cloud computing

The trend toward cloud computing started in the late 1980s with the concept of grid computing when, for the first time, a large number of systems were applied to a single problem, usually scientific in nature and requiring exceptionally high levels of parallel computation.

That said, it's important to distinguish between grid computing and cloud computing. Grid computing specifically refers to leveraging several computers in parallel to solve a particular, individual problem, or to run a specific application. Cloud computing, on the other hand, refers to leveraging multiple resources, including computing resources, to deliver a "service" to the end user.

- *In grid computing, the focus is on moving a workload to the location of the needed computing resources, which are mostly remote and are readily available for use. Usually a grid is a cluster of servers on which a large task could be divided into smaller tasks to run in parallel. From this point of view, a grid could actually be viewed as just one virtual server. Grids also require applications to conform to the grid software interfaces.*
- *In a cloud environment, computing and extended IT and business resources, such as servers, storage, network, applications and processes, can be dynamically shaped or carved out from the underlying hardware infrastructure and made available to a workload. In addition, while a cloud can provision and support a grid, a cloud can also support nongrid environments, such as a three-tier Web architecture running traditional or Web 2.0 applications.*

In the 1990s, the concept of virtualization was expanded beyond virtual servers to higher levels of abstraction—first the virtual platform, including storage and network resources, and subsequently the virtual application, which has no specific underlying infrastructure. Utility computing offered clusters as virtual platforms for computing with a metered business model. More recently software as a service (SaaS) has raised the level of virtualization to the application, with a business model of charging not by the resources consumed but by the value of the application to subscribers.

The concept of cloud computing has evolved from the concepts of grid, utility and SaaS. It is an emerging model through which users can gain access to their applications from anywhere, at any time, through their connected devices. These applications reside in massively scalable data centers where compute resources can be dynamically provisioned and shared to achieve significant economies of scale. Companies can choose to share these resources using public or private clouds, depending on their specific needs. Public clouds expose services to customers, businesses and consumers on the Internet. Private clouds are generally restricted to use within a company behind a firewall and have fewer security exposures as a result.

The strength of a cloud is its infrastructure management, enabled by the maturity and progress of virtualization technology to manage and better utilize the underlying resources through automatic provisioning, re-imaging, workload rebalancing, monitoring, systematic change request handling and a dynamic and automated security and resiliency platform.

### The dynamic data center model

As more and more players across the IT industry announce cloud computing initiatives, many CIOs are asking IBM how they can adopt cloud computing technologies and management techniques to improve the efficiency and flexibility of their own data centers and other computing environments. In response, IBM has recently introduced its vision of data centers which support a dynamic infrastructure that unifies the strengths of the Web-centric cloud computing model and the conventional enterprise data center (as shown in Figure 2).

These data centers will be virtualized, efficiently managed centers, which will employ some of the tools and techniques adopted by Web-centric clouds, generalized for adoption by a broader range of customers and enhanced to support secure transactional workloads. With this highly efficient and shared infrastructure, it becomes possible for companies to respond rapidly to new business needs, to interpret large amounts of information in real time and to make sound business decisions based on moment-in-time insights. The data center that supports a dynamic infrastructure is an evolutionary new model that provides an innovative, efficient and flexible approach in helping to align IT with business goals.

The remaining sections of this paper provide a high-level description of the data center in support of a dynamic infrastructure with underlying technology enablers, such as virtualization, automation, provisioning, monitoring and capacity planning. Finally, examples of actual data center implementations are discussed to reveal which characteristics of dynamic infrastructure models can offer the most value to customers of any size, across a variety of usage scenarios.
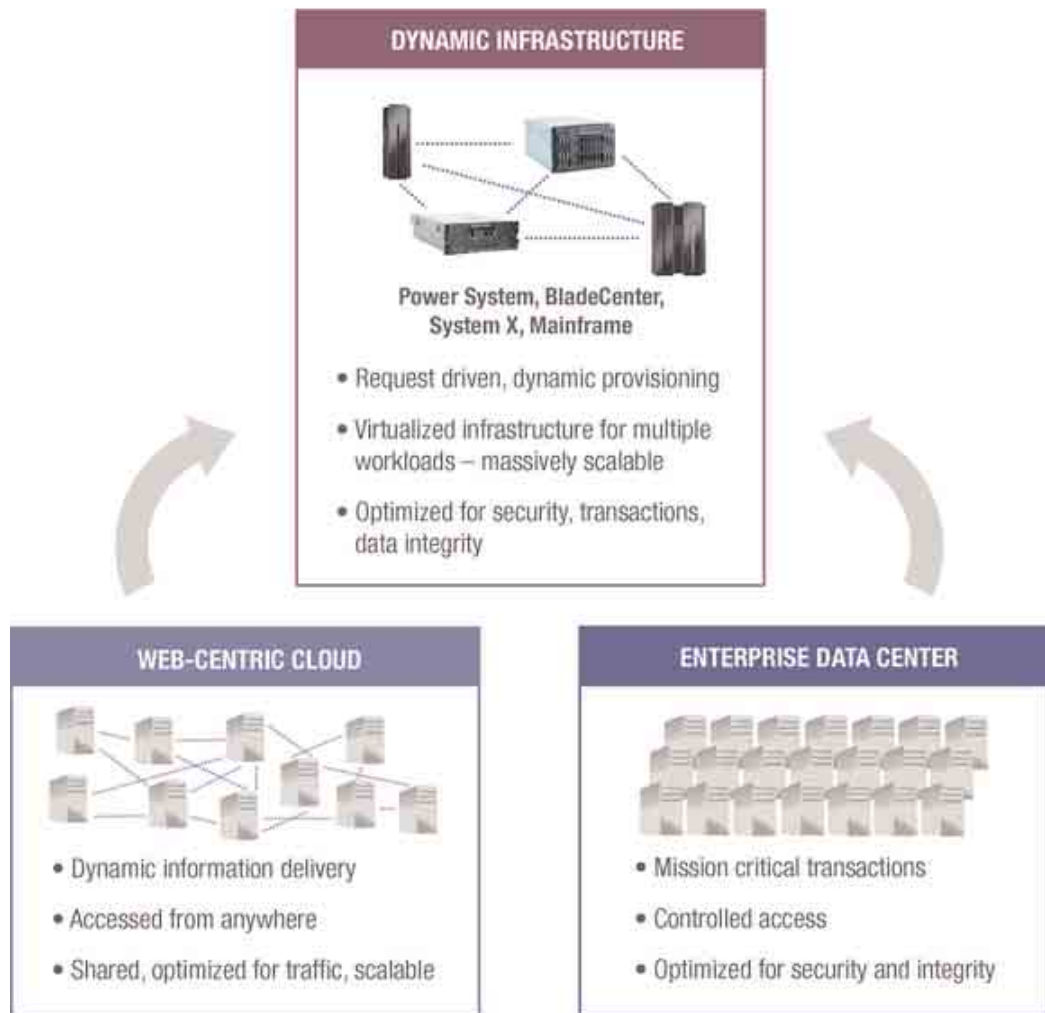
## CLOUD COMPUTING AND THE DYNAMIC ENTERPRISE DATA CENTER



**DYNAMIC INFRASTRUCTURE**

Power System, BladeCenter,
System X, Mainframe

- Request driven, dynamic provisioning
- Virtualized infrastructure for multiple workloads – massively scalable
- Optimized for security, transactions, data integrity

**WEB-CENTRIC CLOUD**

- Dynamic information delivery
- Accessed from anywhere
- Shared, optimized for traffic, scalable

**ENTERPRISE DATA CENTER**

- Mission critical transactions
- Controlled access
- Optimized for security and integrity

*Figure 2. Cloud computing and the data center that supports a dynamic infrastructure*

## Architecture framework and technology enablers

From the high-level architectural point of view, dynamic infrastructure services can be logically divided into layers, as shown in Figure 3. The physical hardware layer is virtualized to provide a flexible, adaptive platform to improve resource utilization. The keys to dynamic infrastructure services are the next two layers: virtualization environment and management. The combination of these two layers ensures that resources in a data center are efficiently managed and can be provisioned, deployed and configured rapidly. In addition, a dynamic infrastructure is designed to handle a mixture of workload patterns, as discussed in the *Business use cases* section of this paper.
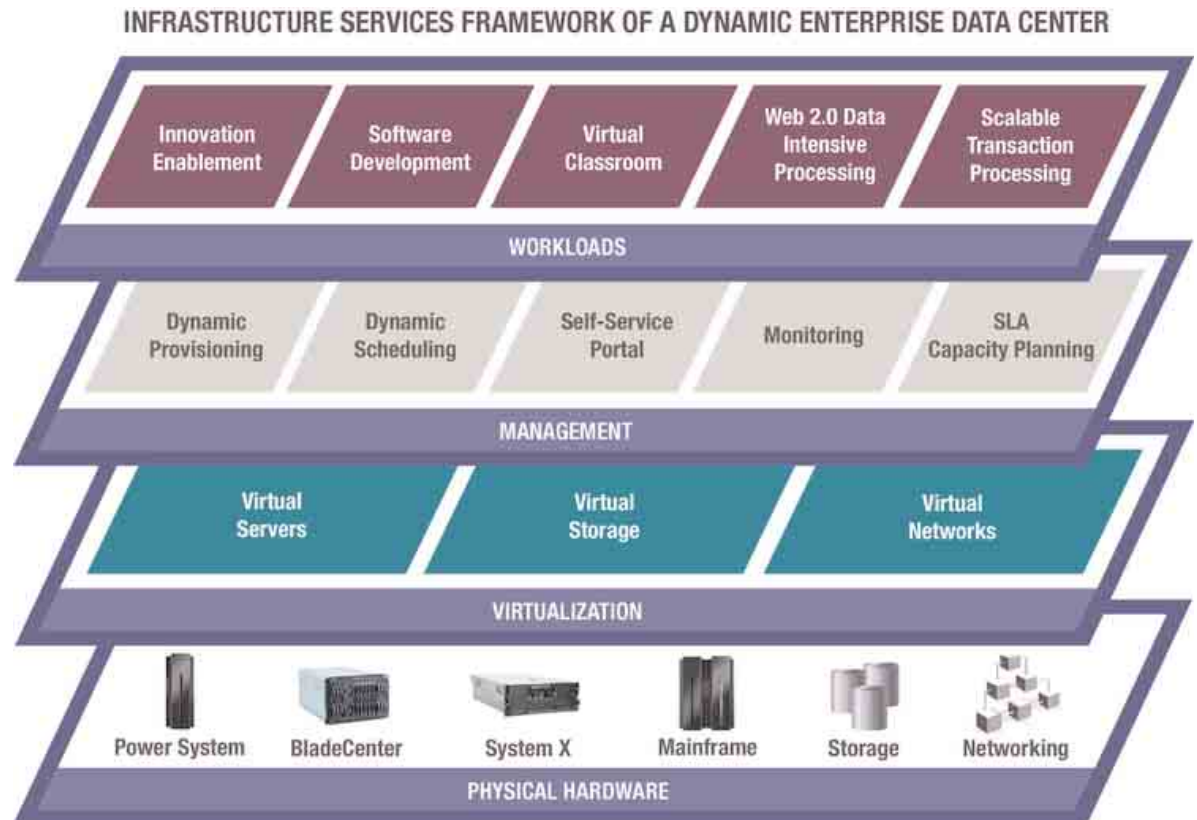


*Figure 3. Infrastructure services framework of a data center that supports dynamic infrastructure*

### Virtualized environment

#### What is virtualization?

Virtualization refers to the abstraction of logical resources away from their underlying physical resources in order to improve agility and flexibility, reduce costs and thus enhance business value. In a virtualized environment, computing environments can be dynamically created, expanded, shrunk or moved as demand varies. Virtualization is therefore extremely well suited to a dynamic cloud infrastructure, because it provides important advantages in sharing, manageability and isolation (that is, multiple users and applications can share physical resources without affecting one another). Virtualization allows a set of underutilized physical servers to be consolidated into a smaller number of more fully utilized physical servers, contributing to significant cost savings.

There are many forms of virtualization commonly in use in today's IT infrastructures, and virtualization can mean different things to different people, depending on the context. A common interpretation of server virtualization is the mapping of a single physical resource to multiple logical representations or partitions. Logical partitions (LPARs) and virtual machines (VMs) are examples of this definition; IBM pioneered this space in the 1960s.

Virtualization technology is not limited to servers, but can also be applied to storage, networking and application, which could be the subject of their own papers.

#### How does server virtualization work?

In most cases, server virtualization is accomplished by the use of a hypervisor to logically assign and separate physical resources. The hypervisor allows a guest operating system, running on the virtual machine, to function as if it were solely in control of the hardware, unaware that other guests are sharing it. Each guest operating system is protected from the others and is thus unaffected by any instability or configuration issues of the others. Today, hypervisors are becoming a ubiquitous virtualization layer on client and server systems. There are two major types of hypervisors: bare-metal and hosted hypervisors.

#### *Bare-metal hypervisors*

A bare-metal hypervisor runs directly on server hardware to provide virtual machines with fine-grained timesharing of resources. Examples of firmware-based bare-metal hypervisors include IBM System z® Processor Resource Systems Manager™ (PR/SM™) and the IBM System i® and IBM System p® POWER® Hypervisor. Examples of software-based bare-metal hypervisor include z/VM®, VMware ESX Server, Microsoft® Hyper-V™ and Xen Hypervisor. The overhead of firmware-based hypervisors is generally less than the overhead of software-based hypervisors. As a result, virtualization implemented at the server hardware level can provide the highest efficiency and performance.

System z PR/SM is a hardware partitioning capability that enables highly scalable, robust, reliable virtual server hosting for multiple operating system images, each in its own LPAR. Virtualization is built in to the System z, which provides the software-based z/VM. Because z/VM virtualization technology enables highly granular sharing of system resources at high levels of efficiency and resource utilization, it allows clients to run hundreds, or even thousands, of Linux® servers on a single mainframe in parallel with System z operating systems.

System p enables users to request virtualized LPARs with IBM AIX® or Linux operating systems. System p has a micropartitioning capability, that allows the system to assign partial CPUs to LPARs. A partial CPU can be as granular as 1/10 of a physical CPU. Thanks to this micropartitioning and virtual I/O capability, the number of LPARs supported by a System p is significantly increased. Furthermore, the LPAR's CPU resource can be managed by IBM Enterprise Workload Manager™, which monitors CPU demand and usage and employs business policies to determine how much CPU resource is assigned to each LPAR. Micropartitioning and virtual I/O make a powerful virtualized infrastructure available for cloud users.

Using technologies such as VMware and Xen, virtualization is now extended to the x86 platforms. With these technologies, many virtual servers with different operating systems can coexist on and share the resources of a single x86-based server. Consolidation of applications running in heterogeneous operating system environments and on underutilized hardware servers becomes possible.

Speedy provisioning is another benefit of using VMware and Xen. Because a virtual server's operating system images and configuration and log files are stored as file sets on the file system of the host or management console server, operating system images of a virtual machine can be run on one physical server and moved or copied transparently to another. Provisioning of a new server or cloning of an existing server can be accomplished by simply creating a virtual server on an existing physical system and copying the previously saved virtual server images. In other words, a server can be provisioned without reinstalling the operating system or the applications running on it.

### *Hosted hypervisors*

A hosted hypervisor runs on a host operating system and uses operating system services to provide timesharing of resources to virtual machines. Examples of software-based hosted hypervisors include VMware Server and Microsoft Virtual Server.

## Infrastructure management

*Virtualization is the fundamental technology enabler for the infrastructure services of the data center. Achieving best business value from virtualization, however, will require a management layer that acts like the brain or control center to efficiently manage the resources in the whole environment. The following sections review the major components of this layer.*

### Automation

Infrastructure administration is one of the major challenges in a virtualization environment. Simply building a virtualization environment without the proper approach to administration can increase complexity and thus generate added costs—costs high enough to cancel out the cost savings derived from virtualization in the first place.

Automation is the key to managing these problems. It is critical that a cloud be equipped with tools that facilitate, simplify and enable management of the physical environment that provides the virtual server resources.

### Automated provisioning

Automation is an important technique that is applied to two of the most frequent tasks performed in a dynamic data center: the onboarding and offboarding of applications. Onboarding is the process of installing and configuring the operating system and additional software on servers so that they can be made available to do useful work. Offboarding refers to the steps necessary to automatically reclaim a server so that it is available for other purposes.

Onboarding an application typically starts with provisioning servers. When done manually, onboarding is a time- and labor-consuming process consisting of many complex steps, such as installation of the operating system and software and configuration of the network and storage. These tasks are often error-prone and typically require highly skilled administrators specializing in the area of system, storage and network. Furthermore, applications typically have unique installation and configuration steps, and here too, the potential for human error is a significant risk. Mitigating that risk is possible through automation, by which the many complex tasks involved in onboarding can be carried out on a completely consistent basis.

IBM Tivoli® Provisioning Manager is a key component of IBM's cloud computing solution. It enables cloud administrators to write workflows that automate the installation and configuration of new servers, middleware and applications, and thus it delivers the speedy and efficient construction and management of IT resources.

### Automated reservations and scheduling

Critical to administering computing resources is the ability to understand what the current and future capacity is to accommodate new requests. Without this understanding, one can neither accurately forecast how many customers can be supported, nor ensure that a steady pipeline of applications can be maintained.

A cloud should be able to communicate the provisioning status and availability of resources, provide the capability to schedule the provisioning and deprovisioning of resources and reserve resources for future use.

### Self-service portal

A self-service portal provides systematic request handling and change management capabilities. This is needed to allow customers or customer representatives to request computing resources (or view status of currently deployed services) through a Web portal. A cloud data center must be able to flexibly handle and execute change requests rapidly to align with fast-changing business needs.

The portal also provides administrators the capability to approve or reject computing resource requests. Once approved, provisioning of computing resources, operating systems and the appropriate application stack is carried out automatically by the system.

A request-driven provisioning system should be employed to take user requests for new services or change requirements for existing services; such a system empowers users to extend the service end date or add or remove resources. Figure 4 shows a sample screenshot for adding or removing resources to a service (project), and Figure 5 shows a sample screenshot for changing the project end date.



Figure 4. Sample screenshot for adding resources to a project



Figure 5. Sample screenshot for changing project end date

Another function of the cloud's self-service portal is to allow users to perform common tasks, such as starting, stopping or restarting the servers. Figure 6 shows these functions. Making these self-service tasks available to the cloud's users removes the burden on IT of locating and scheduling an admin to perform them. User satisfaction climbs and administration costs fall.

*Figure 6. Self-service capability to start, stop and restart virtual machines*

### Monitoring

Monitoring resources and application performance is an important element of any environment. The task of monitoring becomes harder, yet more critical, in a virtualized environment. The benefits provided by monitoring include:

- *Collecting historic data to assist with planning future data center resource needs and to optimize virtualized resource placement;*
- *Capturing real-time data to quickly react to unexpected resource needs;*
- *Measuring adherence to performance service level agreements (SLAs);*
- *Proactively generating alerts and detail data to quickly detect and solve application problems; and*
- *Reporting resource usage data by application, necessary for allocating costs appropriately.*

IBM uses IBM Tivoli Monitoring to monitor the health (CPU, memory and storage) of the servers provisioned by the cloud. This involves installing a monitoring agent on each cloud server and configuring the monitoring server. The agents collect information from cloud resources and periodically transfer that data to the monitoring data warehouse.

Individual servers or collections of servers can be monitored. Detailed information on each monitored resource can be viewed with ITM and fully integrated with the cloud portal. Also, summary information denoting server health can be viewed directly from the cloud self-service portal. Figure 7 shows the CPU, memory and disk summary information that is consolidated at a project level where a project could contain more than one server or resource.



*Figure 7. Projects summary*

### Capacity planning

The cloud computing model reduces the need for capacity planning at an application level. An application can simply request resources from the cloud and obtain them in less than an hour in accordance with dynamic demand. Thus, it is far less important to correctly predict the capacity requirements for an application than it is in traditional data centers, for which as many as six months might be needed to order and install hardware dedicated to the application.

On the other hand, virtualization makes it harder and more important to plan capacity from the data center's perspective. In the past, data center managers could use the projections from applications, take into account the hardware on order, and thus avoid having to dynamically adjust the capacity of deployed hardware. Traditionally, a data center would just need to make sure that it had the capability to support the hardware planned by individual applications. In a cloud environment, however, many different applications will be installed. It becomes the data center manager's responsibility to predict the average or total resource requirement of all the applications and to order enough hardware in advance independently of the input from application owners.

The basis for capacity planning, then, lies in monitoring existing usage and keeping track over historical time periods. Long-term trends can be projected based on previous activity and adjusted without any knowledge of business plans. In a data center-driven cloud, typical capacity planning techniques can be applied for the most part. Since clouds use virtualized resources that share the same physical resources, this makes capacity planning somewhat more complex. In contrast, the capacity planning does not need to consider each individual application, and can simply track and project the overall summation of all applications on the cloud.

There are still times when an individual application may be significant to the overall data center, or require initial capacity without allowing for its capacity to grow organically over a longer time period. To address this need, IBM uses its Performance and Capacity Estimation Service (PACES), previously known as Sonoma. PACES is a Web-based service for evaluating performance and planning the capacity of business patterns, product workloads and custom user-defined workloads. The PACES capacity planning engine is based on a mathematical model involving queuing theory that has been fine-tuned using empirical data from numerous past engagements. IBM regularly validates the underlying model based on the results of recent projects.

Currently, PACES is used by IBM employees around the world from groups such as IBM TechLine group, IBM Sales, IBM Global Services and IBM IT Architects. Customers come to IBM with an application framework and performance objectives. PACES models these complex relationships and either suggests an appropriate configuration solution or projects the performance of a user-specified configuration. Today, PACES supports a large library of different kinds of workloads, along with an extensive library for such hardware as IBM System x®, System p, Sun and HP. Figure 8 shows a sample PACES objective specification screen, and Figure 9 shows a sample PACES estimate results screen.

*Figure 8. Sample PACES objective specification screen*



*Figure 9. Sample PACES estimate results screen*

## Business use cases

With both a long-term pioneering position in the virtualization space and decades of experience in designing and running large, mission-critical distributed computing environments for clients, IBM has recently introduced its vision of a data center that supports dynamic infrastructure.

Leading customers, such as Google and the Government of Wuxi in China, have worked collaboratively with IBM to develop this vision by verifying best practices for running data centers with Web 2.0 and other types of workloads. The following sections share some of the insights that have been gained.

### Innovation enablement

**IBM's internal innovation portal**

The office of the IBM CIO created an internal innovation program to enable and foster a culture of innovation and to harvest the collective intelligence of our large employee base. This program allows internally developed projects to be publicized and made available for use by the IBM internal (and self-selected) *early adopters*, who provide valuable feedback as they use new applications. Innovators not only learn what application features and functions users like, they also receive feedback on nonfunctional areas such as performance, scalability, operational guidelines, resource utilization and support requirements. The innovation program staff then collects metrics and analyzes the business value of the program.

Previously, a typical pilot team needed four to twelve weeks to identify, procure and build a pilot infrastructure, and additional time to build a security-compliant software stack, so that developers could begin building or deploying applications. Subsequently, the CIO office needed to dedicate highly skilled administrators to manage this IT environment. One of its primary business goals was to reduce IT management costs while accelerating innovation.

To achieve this goal, IBM implemented a private cloud to serve as an innovation portal. This cloud was developed via a virtualized infrastructure dedicated to hosting applications and their associated services environment. Innovators who want to obtain broad, user-based feedback can publish information about their applications on this role-based, self-service portal. They then deploy these applications by filling out a form defining the required hardware platform, CPU, memory, storage, operating system and middleware, as well as the project team members and their associated roles. This process takes about five minutes. After submitting the request through the portal, the program administrator is notified and logs in to approve, modify and/or reject the request. If approved, the system begins a process to build the server(s). The process is fully automatic and completes in approximately an hour or less.

This implementation contained all the core elements of a cloud, including a service-oriented architecture with Web services, provisioning, security compliance and monitoring engines, as well as virtualization technologies. The various IBM solutions used to perform these functions include IBM Tivoli Provisioning Manager, IBM Tivoli Security Compliance Manager, IBM Enterprise Workload Manager, Remote Deployment Manager, Cluster Systems Management, Network Installation Manager and IBM Tivoli Monitoring.

This business-responsive data center has proven enormously successful in fostering innovation. In 2007, the innovation portal hosted over 110,000 registered early adopters (almost one third of the IBM employee population). Over 70 different applications were deployed, used and evaluated. Many of these applications graduated out of the program into production use. Ten of these applications were subsequently turned into IBM products. One of the most successful projects was IBM Lotus® Sametime® Version 7.5—Web Conferencing, which had 65,000 registered users providing feedback and comments over blogs and forums offered through the innovation portal.

**Sogeti innovation cloud**

IBM's success in driving innovation through cloud computing has been mirrored by IBM clients as well. One such client is Sogeti, a large consulting company in Europe and a subsidiary of CapGemini, which wanted its managers to convene at their annual meeting and brainstorm new ideas.

Like many other companies, Sogeti was struggling with this question: Where do we start? They had neither the right tool to enable this kind of dialogue among the managers, nor the resources and time to implement a solution on their own. During the progression of ideas through the innovation life cycle, IT resources are often necessary to support the activities. A barrier to this process for Sogeti, as well as many other customers, is the time required to request, instantiate and support a set of IT systems for an innovation program.

Based on experience working with different clients to create an innovation program, IBM believes that collaboration tools alone do not yield desired results as effectively as a structured innovation platform and program. Figure 10 shows the logical architecture of an innovation platform, including each of the *engines* required for the application components and features. Information about the application and associated people is stored in a common IBM Idea Factory database and an LDAP repository, respectively. An optional security authentication and authorization component is shown as well, for enhanced security and customization for granular content authorization.



**LOGICAL ARCHITECTURE OF AN INNOVATION PLATFORM**

*Figure 10. Logical architecture of an innovation platform*

IBM Idea Factory is a specific solution developed by IBM to implement the generic innovation platform. Idea Factory is an Enterprise Web 2.0 solution that accelerates the innovation life cycle. The innovation life cycle is usually divided into five phases: *ideation*, *idea selection*, *development*, *incubation* and *graduation*. The goal of Idea Factory is to provide a flexible, dynamic platform for innovators to collaborate on ideas and create customized Websites for interacting with interested users.

For Sogeti, IBM decided to leverage an IBM Idea Factory solution and host it in a regional IBM cloud data center, based on a cloud data center model. As mentioned above, one barrier that Sogeti had was that during the progression of ideas through the innovation life cycle, IT resources must be requested, instantiated and supported, which requires extensive time and highly skilled IT system administrators.

Cloud computing overcomes this barrier perfectly. IBM, therefore, decided to host this Idea Factory solution in a cloud environment in which VM images representing different components within a standard innovation platform could be deployed easily and quickly as needed. By using this virtualization technique within a cloud environment, IBM was able to not only automate the process and reduce the time and effort needed to provision a new set of J2EE application systems from days or weeks to a matter of minutes or hours, but also make the whole application easier to scale on demand.

In this way, IBM helped Sogeti achieve its goal of empowering its employees to easily share and collaboratively develop innovative ideas and concepts more effectively.

### Software development and test environments
### China Cloud Computing Center at Wuxi

The city of Wuxi, located about 100 miles outside of Shanghai, China, has an economic development project to establish a software park (a specialized zone that provides substantial tax incentives for businesses that open new offices in the software park). One challenge facing startups in Wuxi was the high upfront investment in IT infrastructure they needed to make before being able to accept business from enterprise clients.

To address this challenge and attract companies to the software park, the municipal government of Wuxi worked with IBM to build a cloud computing center based on the dynamic infrastructure model. Tenants in the software park can now use this data center to rent software development and test environments. Figure 11 shows a logical view for its management and customer environments.
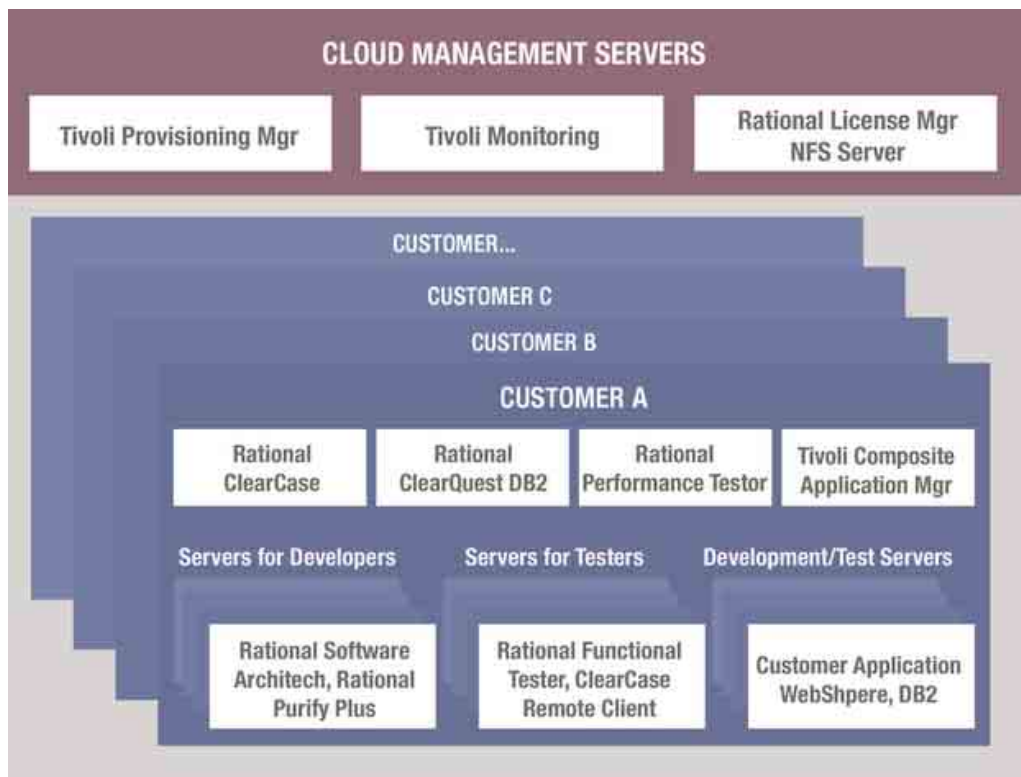


*Figure 11. Management and customer environments*

This cloud utilizes Linux kickstart (a scripting tool for automating Linux OS installation), network file system (NFS), logical volume manager (LVM) and either Xen or the p5 hypervisor as the underlying virtualization environment. The primary products that IBM used in this cloud implementation were IBM Tivoli Provisioning Manager and IBM Tivoli Monitoring. An easy-to-use Web 2.0 user interface and business management logic components have been implemented for users to easily request and manage projects consisting of collections of virtual machines.

Besides automatically provisioning operating systems like Linux Red Hat Version 5 on both System x and System p, and some middleware like WebSphere® Application Server Version 6.1 and DB2® Enterprise Server Edition Version 9.1, this solution also makes it possible to automatically provision many Rational® products to provide software companies with best-in-class development and test environments. IBM Rational solutions utilized for this purpose include IBM Rational for Multiplatform Version 7.0.1 Multilingual, IBM Rational Performance Tester for Mulitplatform Version 7.0 Multilingual, IBM Rational PurifyPlus Enterprise Edition for Mulitplatform V 7.0 Multilingual and IBM Rational Software Architect for Mulitplatform Version 7.0.1 Multilingual.

Because multiple customers are hosted within one environment, this solution requires exceptionally effective network isolation and security. In this virtualized environment, hosts from one physical server may have VMs used for multiple projects; one project might also span multiple hosts. Virtual private network (VPN) technology is used to make sure each client has its own isolated network. When resources are provisioned, additional networks/bridges are configured on either the Xen host or virtual I/O server.

### Advanced computing model for data-intensive workloads
### IBM/Google Academic Initiative

Today's leading Web 2.0 initiatives face special challenges. With more Web 2.0 applications freely available to Internet users, and with more users uploading audio and video, two problems are commonly faced by companies like Google and Facebook: how to reliably store all their data and how to extract the maximum business value from their large volumes of daily Web traffic.

One solution to such problems lies in the MapReduce distributed parallel programming model, which is increasingly used to write these types of business analytics programs. MapReduce allows the processing to be distributed across hundreds to thousands of nodes, all of them working in parallel on a subset of the data. The intermediate results from these nodes are combined, sorted and filtered to remove duplicates before arriving at the final answer.

The majority of processing at Google is based on this MapReduce model. Google's success has prompted other companies to follow the same model, and the increasing popularity of MapReduce-style programming inspired the Apache Hadoop project, which is an open-source implementation of the MapReduce programming framework. Many companies use Apache Hadoop for large-scale business analytics, ranging from understanding users' navigation patterns and trends on their Websites to building targeted advertising campaigns.

The characteristics of the MapReduce programming model require an underlying compute infrastructure that is highly scalable. A data center platform that supports dynamic infrastructure provides an ideal foundation for such workloads.

To promote this surge of interest in MapReduce-style programming, IBM and Google announced a partnership in October 2007 to provide a number of data centers for use by the worldwide academic community. These centers are powered by an IBM solution based on dynamic infrastructure architecture for data center management, which allows users to quickly provision large Hadoop clusters for students who might otherwise be short of required IT resources to complete their lab assignments or run their research programs. Some of the leading universities using these centers and teaching courses on the MapReduce methods are University of Washington, University of Maryland, Massachusetts Institute of Technology, Carnegie Mellon University, University of California Berkeley and Colorado State University.

The heart of this solution is to automatically provision a large cluster of virtual machines for students to access through the Internet to test their parallel programming projects. As a result, physical machines, or virtual machines created using the Xen hypervisor, can be provisioned rapidly and automatically using

Network Installation Manager, Remote Deployment Manager or Cluster Systems Manager, depending upon the operating system and platform. The cluster is powered with open source software, including Linux (Fedora), Xen systems virtualization and the Hadoop workload scheduler. MapReduce Tools for Eclipse, which is open source software designed by IBM to help students develop programs for clusters running Hadoop, is available in Hadoop 0.16.2. Although the current implementations of this cloud support Xen specifically, the framework also allows for other software virtualization technologies such as VMware ESX Server.

From this joint effort with Google, IBM has not only gained a much deeper understanding of the characteristics and programming model for data-intensive workloads such as business analytics applications, but also invaluable hands-on experience concerning how the dynamic infrastructure model can be applied to this type of workload.

## Summary

Today's IT realities make cloud computing a good fit for meeting the needs of both *IT providers* who demand unprecedented flexibility and efficiency, lower costs and complexity and support for varied and huge workloads and *Internet users* who expect exceptionally high availability, function and speed.

As technology—such as virtualization and corresponding management services like automation, monitoring and capacity planning services—become more mature, cloud computing will become more widely used for increasingly diverse and mission-critical workloads.

IBM can help. With its distinguished history of managing and supporting the largest and most complex IT infrastructures, and with its long-term pioneering position in the virtualization space, IBM is well-positioned to implement its vision of the cloud computing data center to meet clients' needs and solve their challenges—today and tomorrow.

## References

HiPODS papers related to or referred to in this paper include:

- *Creating a platform for innovation by leveraging the IBM Idea Factory solution, March 2008 at*
  http://download.boulder.ibm.com/ibmdl/pub/software/dw/wes/hipods/Idea_Factory_wp_14Mar.pdf
- *Sonoma: Web Service for Estimating Capacity and Performance of Service-Oriented Architecture (SOA) Workloads, October 2006 at* www.software.ibm.com/software/dw/wes/hipods/SONOMA_wp9Oct_final.pdf

See all the HiPODS white papers at www.ibm.com/developerworks/websphere/zones/hipods/library.html

Of particular interest, see also:

- *Building a Smarter Planet with a Dynamic Infrastructure, February 2009* http://www.ibm.com/common/ssi/fcgi-bin/ssialias?infotype=SA&subtype=WH&appname=STGE_OI_IS_USEN&htmlfid=OIW03021USEN&attachment=OIW03021USEN.PDF
- *IBM's Vision For The New Enterprise Data Center, December, 2008* http://www.ibm.com/common/ssi/fcgi-bin/ssialias?infotype=SA&subtype=WH&appname=STGE_OI_IS_USEN&htmlfid=OIW03013USEN&attachment=OIW03013USEN.PDF

Here are links to other related references in this paper:

- *HousingMaps –* www.housingmaps.com
- *FlickrVision –* flickrvision.com
- *VMware Vmotion –* vmware.com/products/vi/vc/vmotion.html
- *IBM Tivoli Provisioning Manager –* www.ibm.com/software/tivoli/products/prov-mgr
- *IBM System Director Active Energy Manager –* www.ibm.com/systems/management/director/extensions/actengmrg.html
- *Apache Hadoop –* hadoop.apache.org
- *MapReduce – MapReduce: Simplified Data Processing on Large Clusters at* labs.google.com/papers/mapreduce.html
- *Green Computing Initiatives –* www.ibm.com/technology/greeninnovations/index.html
- *IBM-Google academic initiative –* www.ibm.com/jct09002c/university/scholars/skills/internet-scale/index.html

Recyclable, please recycle

OIW03022-USEN-00