

# **Cloud Management**

Nigel Cook, Dejan Milojicic, Vanish Talwar

HP Laboratories HPL-2011-238

# Keyword(s):

cloud services; service management; middleware; heterogeneity; integration; scalability; service level agreements

# Abstract:

Cloud Computing offers a number of benefits, such as elasticity with the perception of unlimited resources, self-service, on-demand, automation, etc. However, these benefits create new requirements for management of Cloud computing. On the back-end, economic limitations dictate careful consolidation of servers with clear sustainability analysis; managed levels of abstractions are higher (from hardware, to VMs, to services); and reliability, availability, and supportability are built into higher levels of systems and services. On the client-side, Cloud services have to be easy to use/manage, perform well, and be reliable. On both sides, geographical distribution and its implications on business continuity is a rule rather than exception; scalability is built-in by design; and QoS is still being defined. In this paper, we discuss new requirements and approaches to Cloud management. We present a few examples of Cloud management for private, public, and HPC Clouds. Based on these, we derive conclusions about manageability of current platforms and then make predictions about the research challenges of future Cloud management. We expect these findings to help designers of next generation hardware and software platforms to develop more manageable systems and solutions.

External Posting Date: December 8, 2011 [Fulltext]

Internal Posting Date: December 8, 2011 [Fulltext]

Approved for External Publication

# **Cloud Management**

Nigel Cook • Dejan Milojicic • Vanish Talwar

## Abstract

<sup>1</sup>Cloud Computing offers a number of benefits, such as elasticity with the perception of unlimited resources, selfservice, on-demand, automation, etc. However, these benefits create new requirements for management of Cloud computing. On the back-end, economic limitations dictate careful consolidation of servers with clear sustainability analysis; managed levels of abstractions are higher (from hardware, to VMs, to services); and reliability, availability, and supportability are built into higher levels of systems and services. On the client-side, Cloud services have to be easy to use/manage, perform well, and be reliable. On both sides, geographical distribution and its implications on business continuity is a rule rather than exception; scalability is built-in by design; and QoS is still being defined. In this paper, we discuss new requirements and approaches to Cloud management. We present a few examples of Cloud management for private, public, and HPC Clouds. Based on these, we derive conclusions about manageability of current platforms and then make predictions about the research challenges of future Cloud management. We expect these findings to help designers of next generation hardware and software platforms to develop more manageable systems and solutions.

*Keywords: cloud services, service management, middleware, heterogeneity, integration, scalability, service level agreements.* 

Abbreviations: QoS: Quality of Service; SLA: service level agreements; IT: Information Technology; DevOps: Development Operations; NVRAM: Nonvolatile Random Access Memory; AWS: Amazon Web Services; VM: virtual machines; CAPEX/OPEX: Capital/Operational Expenditure; SSD: Solid State Disks; WBEM: Web-Based Enterprise Management.

#### 1 Introduction

Cloud computing is an emerging paradigm, with growing popularity and adoption [1]. Cloud providers host

N. Cook Hewlett Packard, Littleton, Colorado, USA nigel.cook@hp.com

D. Milojicic Hewlett Packard, Laboratories, Palo Alto, CA, USA dejan.milojicic@hp.com

V. Talwar Hewlett Packard Laboratories, Palo Alto, CA, USA vanish.talwar@hp.com shared servers, and deliver computing, storage, and software to end-consumers as a service. Both Gartner and IDC have estimated healthy growth of cloud computing adoption [2][3]. Cloud services include compute-ondemand, online storage, online/shared office applications, key value store, and email, among many others services. Examples of public cloud providers are Amazon AWS [4], GoGrid [5], and Rackspace [6]. Several other companies have cloud offerings, such as HP [7], Google [8], IBM [9], and Microsoft [10].

Traditional Web companies, such as Google and Yahoo have proprietary Cloud management stacks. Amazon was among the first to publish their interfaces for Cloud, including management. Eucalyptus is an open source implementation of Amazon interfaces [11]. RightScale [12] focuses primarily on Cloud management aspects of Clouds. Most recently, OpenStack [13] is an effort to develop a Cloud stack by a number of companies (over 130 at the time of writing this paper and growing). In addition, there are other open source Cloud stack efforts under way, such as OpenNebula [14] and Tashi [15]. Research efforts and testbeds include RESERVOIR [16], Open Cirrus [17], and Open Cloud Consortium [18]. Other examples of Cloud management among many include CloudWatch, Nimsoft, MMC, Mesos [19], Monalytics [20][21], vManage [22], and multiple managers [23].

Traditional standardization organizations, such as *DMTF*, *NIST*, and *IEEE*, have independent efforts in standardizing different aspects of Clouds and Cloud management. They are still early in the process to understand the impact of these efforts. Amazon Web Services interfaces appear to be a de facto standard interface, while OpenStack is getting momentum as an open source implementation thereof.

Cloud computing is enabled by advances in virtualization, service oriented computing, and utility computing. There are several requirements for cloud computing to be successful. These include low-cost, SLA compliance, security guarantees, high availability, energy efficiency, and accurate accounting. The key to meeting these requirements is effective management of cloud resources and services. This covers all aspects of the data center lifecycle from bring-up, provisioning, scheduling, monitoring, failure management, and shutdown.

As IT becomes increasingly automated, so does the importance of IT manageability. This is especially true in Cloud where automation is essential for driving down the cost. Manageability is defined as the collective processes of deployment, configuration, optimization, and administration during the lifecycle of IT systems and services. Recent examples of Amazon and VMware outages, which impacted the business continuity of a number of hosted companies, are key indicators of the importance of manageability.

Manageability has multiple dimensions. Resource management concerns scheduling and resource assignment, performance and availability, virtual machines, workload, and OS functions. Automation addresses deployment, provisioning, monitoring, configuration, changes, and problems.

Manageability targets managed objects, which can be hardware or software (object, service, data, etc.). The lifecycle of a managed object is presented in Figure 1, from bring-up, through operation, over failures/changes, till retire/ shutdown. A managed object can have different granularity and composition. The lifecycle of a managed object can also be of different duration; in Clouds, it is typically shorter compared to non-Clouds.



Figure 1. Lifecycle of a Managed Object.

While the above figure is true in cases when the full system is owned and managed by the service provider, in the case of Clouds this is not true. Different parts of the system can be managed by different owners and in different domains, behind different firewalls (See Figure 2, red arrows indicating independent management domains).



Figure 2. Managing Clouds and Cloud Services.

Figure 3 shows complexity of different phases and levels of management and how these phases and levels interact. Cloud services are managed at the top of this spectrum, but their management depends on managing objects lower in the dependency chain. Because different objects are managed independently, there is a need for integration of individual managers to avoid inconsistency or undesired behavior.

A distinct feature of Cloud service management is "self-service", typically accomplished through a portal (see Figure 4). An important interplay exists between development and delivery of services. The cloud management environment sits on top of the stack of different layers of Cloud delivery engines, automation engines, and deployment templates and best practices.

Many of the insights in this paper, we based on our prior work in management of Clouds [24], scalable monitoring and analysis [25][26], distributed systems [27], service compatibility [28], SLA management [29], adaptation [30], service deployment [31], federation [32], policy management [33], model-based management [34], change management [35], sustainability [36] and supportability [37]. We have also derived a lot of insights from similar "Future of Software Engineering" workshops, as well as from the specific paper on the future of Middleware [38].



The rest of the paper is organized in the following manner. In section 2, we present three examples of contemporary Cloud management. Section 3 summarizes some of the IT industry trends. In Section 4, we discuss requirements and research challenges. Finally, we summarize the paper in Section 5.



Figure 4. Self Service at the top of Service Management.

#### 2 State of the Art Cloud Management Examples

#### 2.1 Managing Private Clouds: CloudSystem Matrix

HP CloudSystem is an example of a layered management stack for private, but also public or hybrid cloud environments. The environment is constructed as a layering of abstraction as follows:

Virtualization management. The lowest layer provides a lifecycle management of a set of virtualized resources that are drawn from a pool of capacity in the data center. Examples of the virtualized resources and the corresponding management include virtualized servers, storage, and networking, which can be managed by VMware vCenter, Microsoft System Center Virtual Machine Manager, or HP Insight Control. This can be applied to private cloud environments and for public cloud environments, such as OpenStack or Amazon EC2. Each environment provides a notion of an underlying resource capacity implied by a combination of the physical resource being virtualized, or by quotas applied to consumption of individuals or groups. These systems provide capabilities to manage the lifecycle of their virtualized resources, as well as provide monitoring information about the resource consumption and availability of their specific components.

Cloud Service Composition. Built on top of the virtualization management, is the component that manages composition of the virtualization environment to create aggregate cloud service infrastructure. An aggregate service is one that uses a heterogeneous mix of virtualized resources or resource geographies to realize a service offering. Composition of services requires a model of the service components and their relationships, as well as modeling of the capacity and relationships of the underlying virtualized resources. The composition layer uses these two models to schedule use of the virtualized resources to match the infrastructure demand generated by the composite service. Scheduling algorithms take account of service quality considerations, which include both availability considerations, as well as isolation or compliance requirements between different services. This layer monitors the state of the infrastructure elements, alerting on failures, and monitors resource consumption with a goal of providing optimal utilization of the underlying resources, including energy and network bandwidth.

Application Management models the components of a business application and the relationship to the infrastructure provided from the cloud service composition layer. The infrastructure needs of the application can vary by the stage in the lifecycle, or due to varying workload demands placed on the service. As an example, during the development phase of an application, the application may reside on virtualized resources entirely contained within a testbed constructed from public cloud resources, while during production that same application may reside both on an internal private cloud holding the application transaction engine as well as one or more external clouds providing the Web interface and catalog components. While the application is running, the service responsiveness is monitored, and if it falls outside of set limits, then scaling adjustments are made, both by adjusting the number of running application instances, or by requesting adjustment of the infrastructure supplied by the cloud service composition layer.

*QoS Management.* In addition to the DevOps environment (See Section 2.2), there is also a layering of delivery management for applications, which includes scaling of the instances of the application to achieve necessary service levels, maintaining operation in the presence of maintenance cycles, and optimization of facilities utilization by removing unneeded capacity from a service automatically. In order to achieve application management, the application needs to conform to patterns supported by the cloud PaaS layer. The result of this conformance is that the PaaS platform manages the scalability and availability aspects of the services, rather than each application development team needing to create and operate a separate strategy for these aspects.

*Challenges* for enterprise Clouds at the composition layer include algorithms for distributed placement and scheduling of virtualized resources into the distributed capacity pools, particularly for requests targeted at times in the future. For the application management and scaling, a key issue is understanding the scaling model of an application, and interpreting the root cause of application service level changes. Other challenges specific to private, public, and hybrid Clouds include:

- Automated elasticity and SLA guarantees, security, and availability in shared environments is hard to support.
- Unified and integrated management across compute, storage, and network does not exist, preventing end-toend management of applications and Cloud services.
- *Federated management across clouds instances* are hard to achieve for independently managed private Clouds.

#### 2.2 Managing Public Clouds: Internet Data Centers

There has been a recent surge in new Internet companies such as Facebook, Twitter, Google, Amazon, and LinkedIn. These companies provide online services such as search, social computing, and shopping, and they are hosted within large scale and globally deployed data centers accessed by millions of customers/users worldwide. Systems management in such large scale infrastructures provides several challenges. Below we highlight three trends that provide specific challenges and opportunities towards next generation systems management in such infrastructures.

Massive scale in terms of users, machines, data: Existing Internet data centers already contain several hundreds of thousands of machines and this number is increasing to meet the growth in the number of users accessing the online services. A simple back of the envelope calculation easily shows that we can expect several millions of managed objects in such future data centers. This poses several challenges for the automated deployment of OS/VM/application images, load balancing to meet demands, fail-over/reliability of machines and software, as well as capacity planning to ensure service demands are met. Constraints to meet CAPEX, OPEX, and sustainability goals along with requirements to meet guaranteed service levels pose challenges for the design of scalable management systems. Furthermore, large scale systems pose challenges for system logging, monitoring, and analysis for abnormal system behavior to meet high traffic rates. Various frameworks such as Scribe are in use by these companies but they are challenged by increasing scale. The growth of data and its storage poses additional challenges to ensure appropriate dynamic partitioning, migration, and replication to meet service demands, as well as to perform traditional archiving and backup.

Services built by integrating multiple open source frameworks: Internet companies are challenged with the need to reduce the time to bring new services to the market and at the same time ensure scalability. Recent trends include leveraging open source frameworks to quickly bring-up the backend infrastructure in operation at low-cost and leveraging resources to provide better core services. This has resulted in many open source frameworks such as Hadoop, Cassandra, Thrift, Storm, Hive, HBase, MySOL, PHP, Flume etc. The Internet companies integrate these various open source frameworks on Linux to provide their backend and processing infrastructure. While this speeds up the time it takes to bring up the infrastructure, it poses challenges for ongoing operations several and management. First, automated configuration management across multiple tools is a challenge. Gluing together multiple pieces written by different developers requires painful and careful integration and setting of the configuration parameters. Given different possible combinations of the integration, the current processes for configuration are ad-hoc. Further, there are challenges for tuning the framework, both individual and integrated, endto-end.

Furthermore, there are also challenges for end-to-end diagnosis of these integrated frameworks, especially in scenarios where they are pipelined together, e.g. for streaming data processing. Each framework supports a self-managing capability that allows it to recover from failures and abnormalities. However, when these frameworks are integrated together, there is a lack of an end-to-end self-managing capability, and allowing individual self-management loops to proceed without coordination leads to unpredictable behavior and inefficiencies. There is a need to develop an end-to-end monitoring and analysis

framework that can be deployed on-demand in such multistage frameworks.

*DevOps.* A new DevOps model is emerging, i.e. developer and sys-admin operations are merging: several of today's Internet companies develop in-house services and the operations work is also done by in-house system administrators. This implies a culture where development and operations work together with shared responsibility. This is in contrast to previous models where software used to be packaged and shipped. System administrators, who were completely disconnected from the original developers, would deploy the package. An update or new release would occur about once in a year.

In today's Internet companies, releases happen more frequently and do not require physical packaging. Releases take place sometimes weekly or even daily. Agile development methodologies are in use for this new DevOps model. This changes the way administrators and system management tools are designed for deployment and release. Given the shared responsibility, the gap between the silos of program development and operations/admin tasks is disappearing. This implies there is tighter integration between programs and system admin tasks and greater importance for operational efficiency during development. This poses a new model for system management and a new set of tools for this integrated DevOps model. DevOps focuses on application lifecycle management for developers, not end-users, taking products through lifecycle stages: from *package* (application model) through *publish* (environment-specific deployment models); provision and deploy; workload management; and back to package (complete cycle). Specific DevOps functions include:

- Modeling & Configuration Management
- Infrastructure Provisioning
- Application Deployment
- Infrastructure and application monitoring
- Embedded workload management

*Challenges* in this use case include the following:

- *Heterogeneity* of deployment environments, e.g. multiple infrastructure choices, databases, or hypervisors as well as working across private and public Clouds.
- Automated release and testing, to enable stable products (as the versions of managed objects change and the deployed base grows substantially).
- Support and documentation, to resolve issues in a production environment with performance lifecycle management; enough information needs to be captured to enable support to identify problems and provide feedback through DevOps to developers to diagnose and fix issues.
- *Modeling for automated configuration management,* to address complex configurations of service compositions.

- *Maintaining stringent service level guarantees:* to ensure continuous availability of global Internet services with low latency response time even in the presence of flash crowds.
- 2.3 Managing HPC in the Clouds: Towards Exascale

Today's use of Clouds for high performance computing is growing, but it is limited to small scale, testing and development. Amazon has built a top-500 supercomputer in its cloud with 7k cores and achieved speeds of 41.82 teraflops, making it the 231st fastest supercomputer in the world (at the time). They accomplished it with Linux on Intel Xeon X5570 with a 10 Gig Ethernet interconnect. It was de-provisioned soon after running the test but it demonstrated supercomputer-based processing at the price of \$1.60/node hour.

At the high end of HPC, the US Department of Energy is preparing an exascale program, and so are governments of other countries, such as in Europe, China and Japan. An excerpt from the current DOE proposal for exascale computing list of requirements is included below:

|                  | 2010    | 2015   | 2018    |
|------------------|---------|--------|---------|
| Power            | 6MW     | 15MW   | 20MW    |
| Nodes #          | 18,700  | 5,000  | 100,000 |
| Node concurrency | 12      | ~1,000 | ~10,000 |
| Interconnect BW  | 1.5GB/s | 1TB/s  | 2TB/s   |
| MTTI             | Day     | ~Day   | ~Day    |

Table 1. HPC Evolution.

These parameters represent the boundaries of high end HPC, but in many ways they are evolving in a similar direction as high end data centers. The major difference is slower interconnects and less powerful computation nodes, similarity is in power, cooling, and packaging.

In the future, Clouds will contain improved interconnects, such as photonics, that will enable more HPC applications to be executed in the Cloud. The requirements for next generation supercomputers are becoming very similar to Cloud requirements even though some of the design choices may be different.





App's Latency Sensitivity

Figure 5. HPC Applications and Target Platforms.

Of particular interest is differentiating which applications are best suited to which platform. Figure 5 shows the types of applications that best suit Clouds and supercomputers. Applications that exhibit less latency sensitivity and can be allocated to 'lower cost' resources are best suited for Clouds.

A management platform that can perform such matching automatically will benefit HPC Cloud adoption.

The following challenges remain for wider adoption of HPC in Clouds:

- *Latency:* current interconnects deployed in Cloud data centers do not offer sufficient performance for HPC applications. Photonics offers some promise for the future.
- *Cost:* to enable Clouds for HPC, managing cost and pricing is essential. Existing pricing models will have to be expanded, including physical clusters, job submissions, and future reservations.
- *Power:* as HPC grows in performance, power will continue to be one of the main obstacles both for HPC and HPC in the Cloud. Carefully managing power consumption is critical for reducing power cost (power capping, server consolidation, migration, etc.).
- *Virtualization:* while overheads are of less concern for Cloud applications they limit virtualization use for HPC applications. For example, in HPC applications I/O virtualization is not used at all.
- Security: it will be unacceptable to execute some applications globally due to national security concerns. In addition privacy and export rules limit use to specific regions. Automated management of regulatory compliance will be a key differentiator.

#### **3** IT Industry Trends

New technology development always results in faster, bigger, more reliable devices, such as memory, CPU, interconnect, networks, etc. However, today we are at a point where some new technology transitions will have a lasting impact on management.

*NVRAM* systems will have persistency and low latency storage access, driving the need for low-latency and lightweight management stacks. This will require new management models (e.g., new WBEM) and new hardware monitoring and other management tools.

*Novel memory hierarchies*, multi-core, photonics and advances in networking will change systems design and implementation. Management stacks will need to be optimized, lightweight, and decentralized.

*Power and cooling* dominate OPEX/CAPEX. To limit these costs, interfaces will have to be exposed for system and application power management.

As a result, *operating systems* will get redesigned with built-in management in various components (similar to

SMART in disks). There will be multiple components in the architecture that will contribute to management. Therefore integration and federation of management domains will become important.

*Data-intensive computation* and continuous production of data (from sensors and many other devices) will require the ability to archive, and manage the data lifecycle. Data elasticity is not the same as computation elasticity (stateful v. stateless; continuously produced and updated). Management will have to be intertwined with functional support; boundaries between functional support and management are disappearing.

*New application models* such as social networking and big data will require new management architectures and algorithms. This will result in new management models, which will be application-driven.

## 4 Future of Cloud Management: Requirements and Research Challenges

In this section we summarize some of the requirements and challenges of future Cloud Management.

- 4.1 Future Cloud Management Requirements
- *Global scale* (7-8B users), mobile access by most users, elasticity at this scale.
- *Ease of use* resulting in *short time-to-manage*, using visual tools, analytics, what-if analysis, predictions, etc.
- *Cost efficiency*, understanding the costs of hosting services (infrastructure, services, and business objectives).
- *Support for SLAs with multiple objectives*, ability to make tradeoffs in an easy and predictable way.
- Availability and business continuity. Managing replication at the resources level and at the service level; trading off replication cost for the degree of availability.
- Automated regulatory compliance. Due to the global nature of cloud computing, export and privacy rules need to be verified automatically.
- 4.2 Future Cloud Management Research Challenges

Meeting the above requirements, will expose new research challenges to Cloud management. New challenges are derived from the level of scale, resource limitations (power in particular), reliability at such scale, and complexity of managing data, QoS, and integration. These challenges are discussed in more detail below and also summarized in Table 2 for Cloud management today and for research direction.

• *Management at scale*, global and mobile access will result in unpredictable scale up and down. Elasticity of access also results in elasticity of management.

*Federation* will be a way to address scalability and to connect independently managed Clouds.

- *Sustainability.* Environmental awareness is becoming increasingly regulated and it will become a requirement, not just a desirable feature. Power limitations will drive cost and scale as data centers continue to grow.
- *Reliability and Support:* As scale continues to grow, failure rates will also increase, leaving no choice but to automate support. Support will also move away from reactive towards deferred and proactive. Supportability and reliability will be built into the design across all layers.
- *QoS: SLA management* was always hard and it will only grow in complexity with global access, a wide variety of standard and non-standard interfaces, and different APIs for SLA management. Multiple objectives will result in further complexity.
- *Data management.* With continuous generation of new data from sensors, multimedia data formats, and many other sources, the ability to manage this data, and compress, deduplicate, archive, and dispose of it, according to regulatory compliance, will be a huge challenge.
- *Integration* of management components, and run time *composition*. Increasingly more integrated services will result in even higher complexity of versioning, compatibility, and coordination among multiple management components.
- *Quantifying Cloud Manageability* is a research challenge. Some of the ways to quantify manageability are listed below, but new models and metrics need to be devised:
  - ° Checklist of manageability functions
  - ° Number of steps to manage towards desired state
  - <sup>°</sup> Time to manage (including time to insight)
  - ° Documentability (e.g. lines of management code)
  - ° Elasticity of management (manage at scale)
  - ° Availability and continuity of management
  - ° Ease of use (GUIs, visualization, analytics, etc.)

#### 5 Summary

In this paper, we evaluated Cloud management today and some of the trends that we see coming in the future. We presented three examples of Cloud management: public, private, and HPC. For each, we emphasized challenges for the future of Cloud management. We then related Cloud management trends to the general trends in the IT industry. Based on these trends, we summarized some of the requirements and research challenges of future Cloud management.

Cloud Computing has a fundamental role in the future of society, as most IT is migrating towards the Cloud. As mobile services find their way into the Cloud, it will become even more ubiquitous. The role of Cloud management will become essential – particularly in regard to how scale, DevOps, and QoS are addressed. With the tremendous amount of data expected to be generated, dataintensive operations will become dominant compared to those that are compute-intensive, while sustainability and support will change in the future.

The landscape of the Cloud – at different levels of the stack (hardware, services), as well as roles (developers, operators, users) - will differ substantially from the one today. At the hardware layer new technologies will enable greater scale, requiring increased automation and new reliability techniques. Operating these types of evolving Clouds and their services will require frequent updating, an understanding of business trends, and the ability to perform what-if-analysis. Development of new services will increasingly be the result of the composition with continuous rollouts. Most cloud users will be mobile, and many new users will be from developing countries; these powerful user segments will drive innovation and Cloud services pricing models--and therefore Cloud management. (See also Table 3.) Cloud management is fertile ground for fundamental research in systems, applications, and services.

#### Acknowledgements

We would like to thank organizers of the FOME event, Gordon Blaire and Valerie Issarny for their leadership in forming this event, as well for reviewing our paper. Their feedback significantly improved the content and the form. We are also thankful to Fabio Kon for continued guidance. Figure 3 was derived from the similar figure originally created by Sue Charles. Figure 5 was created by Abhishek Gupta as a part of his internship at HP Labs. Sue Charles and Puneet Sharma reviewed submitted document and provided valuable feedback.

### References

- M. Armbrust et al., "Above the Clouds: A Berkeley View of Cloud Computing," tech. report UCB/EECS-2009-28, 2009.
- [2] Gartner. http://www.gartner.com/it/page.jsp?id=1389313; http://www.gartner.com/it/page.jsp?id=1454221.
- [3] IDC,
- http://www.idc.com/research/cloudcomputing/index.jsp
- [4] Amazon AWS, http://aws.amazon.com/.
- [5] GoGrid, http://www.gogrid.com/.
- [6] RackSpace, http://www.rackspacecloud.com/.
- [7] HP Cloud, http://www8.hp.com/us/en/solutions/solutionsdetail.html?compURI=tcm:245-300983&pageTitle=cloud.
- [8] Google Apps, http://www.google.com/apps/intl/en/business/index.html.
- [9] IBM Cloud, http://www.ibm.com/ibm/cloud/.
- [10] Microsoft Cloud, http://www.microsoft.com/en-us/cloud/.
- [11] Eucalyptus, http://www.eucalyptus.com/.
- [12] RightScale, <u>http://www.rightscale.com/</u>.
- [13] Open Stack, http://www.openstack.org/.
- [14] R. Moreno-Vozmediano, R. S. Montero, and I. M. Llorente, "Elastic management of cluster-based services in the cloud", ACDC '09.

- [15] M. Kozuch, et al., Tashi: Location-aware Cluster Management," ACDC'09 June 2009, Barcelona Spain.
- [16] RESERVOIR, www.reservoir-fp7.eu.
- [17] A., Avetisyan, et al., "Open Cirrus A Global Cloud Computing Testbed," IEEE Computer, vol 43, no 4, pp 42-50, April 2010.
- [18] Open Cloud Consortium http://www.opencloudconsortium.org/.
- [19] Hindman, B., et al., Mesos: A Platform for Fine-Grained Resource Sharing in the Data Center, NSDI 2011, March'11.
- [20] Wang, C., et al., "A Flexible Architecture Integrating Monitoring and Analytics for Managing Large-Scale Data Centers", ICAC 2011.
- [21] Kutare, M., et al, "Monalytics: Online Monitoring and Analytics for Managing Large Scale Data Centers", ICAC 2010.
- [22] Kumar, S., et al., vManage: Loosely Coupled Platform and Virtualization Management in Data Centers, In the proceedings of 6th ICAC, Barcelona, Spain, June 2009.
- [23] Kephart, J., et al, "Coordinating Multiple Autonomic Managers to Achieve Specified Power-Performance Tradeoffs," Proc. of the 4<sup>th</sup> ICAC, IEEE CS, 2007.
- [24] Cook, N., Milojicic, D., Talwar, V., "Managing the Cloud Infrastructure", Book Chapter, In Migrating to the Cloud: For Developers, and Technologists, Elsevier Publishers, 2011 (to appear).
- [25] Wang, C., et al., "Statistical Techniques for Online Anomaly Detection in Data Centers", IM 2011.
- [26] Viswanathan, K. et al., "Ranking Anomalies in Data Centers", NOMS 2012, to appear.
- [27] Adams, R., Brett, P., Iyer, S., Milojicic, D., Rafaeli, S., Talwar, V., "Scalable Management", Book Chapter, In Autonomic Computing: Concepts, Infrastructure, and Applications, CRC Press 2006.
- [28] Becker, K., Pruyne, J., Singhal, S., Lopes, A., Milojicic, D., "Automatic Determination of Compatibility in Evolving Services," International Journal of Web Services Research, 2010, 8(1): 21-40.
- [29] Chen, Y., Iyer, S., Liu, X., Milojicic, D., Sahai, A., "SLA Decomposition: Translating Service Level Objectives into system level thresholds," Journal of Cluster Computing, vol 11, no 3, pp 299-311, September 2008.
- [30] Vambenepe, W., Thompson, C., Talwar, V., Rafaeli, S., Murray, B., Milojicic, D., Iyer, S., Farkas, K., Arlitt, M., "Dealing with Scale and Adaptation of Global Web Services Management," Journal of Web Services Research. 2008.
- [31] Talwar, V., Milojicic, D., Wu, Q., Pu, C., Yan, W., and Jung, G., "Approaches for Service Deployment," IEEE Internet Computing, pp. 70-80, vol. 9, no 2., Mar-Apr 2005.
- [32] Bardhan, S., Hidangmayum, R., McGeer, R., Milojicic, D., RN, V., Feldhaus, F., Roeblitz, T., Yahayapour, R., "Practical Federations," Proceedings of the Fifth Open Cirrus Summit, Moscow, IEEE co-sponsored, June 2011.
- [33] Cai, Z., Chen, Y., Kumar, V., Milojicic, D., and Schwan, K., "Automated Availability Management Driven by Business Policies," Proc. of the 10th IFIP/IEEE Symposium on Integrated Network Mgmt, IM'07, Munich, pp 264-273.
- [34] Rivaldo, R., Chen, Y., Milojicic, D. and Adams, R., "SML Model-based Management," Proceedings of the 10th IFIP/IEEE Symposium on Integrated Network Management (IM 2007), Munich, pp 761-764.

- [35] Shankar, C., et al., "Specification-Enhanced Policies for Automated Management of Changes in IT Systems," Proc. of 20th USENIX LISA'06.
- [36] Bash, C., et al. "Cloud Sustainability Dashboard, Dynamically Assessing Sustainability of Data Centers and Clouds," Proceedings of the Fifth Open Cirrus Summit, Moscow, IEEE co-sponsored, June 2011.
- [37] Connelly, C., et al., "Reiki: Serviceability Architecture and Approach for Reduction and Management of Product

Service Incidents," Proc. IEEE ICWS, pp 775-782, Jul 2009.

[38] Issarny, V., et al., "Service-Oriented Middleware for the Future Internet: State of the Art and Research Directions," JISA - Journal of Internet Application and Services, to appear.

| Management Functionality             | State of the Art                                                                                                                                          | Research Direction                                                                                                                                                                                 |
|--------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Management at Scale and Federation   | Hundreds of thousands of nodes in data centers;<br>zones and service-level integration, incremental<br>scalability; simple visualization.                 | Hierarchies of domains, federations of independently<br>managed data centers and Clouds; visualization analytics<br>at full scale.                                                                 |
| Sustainability                       | Tracking power, $CO_2$ and water usage,<br>minimizing environmental impact; introduction of<br>end-to-end sustainability.                                 | Trading off sustainability for QoS, automated<br>sustainability and SLA management, accounting for<br>sustainability of mobile services delivery.                                                  |
| Support and reliability              | Reactive at the high end with field engineers,<br>deferred at the low end with minimal human use;<br>semi-automated.                                      | Preventive, substantially automated, self-healing and<br>rejuvenation of components; field engineers only used at<br>the very high end.                                                            |
| QoS: SLA management                  | Simple services level objectives.<br>Lack of compliance and enforcing SLAs.<br>No integration with business models.                                       | Multi-objectives, business objectives (pricing, costing).<br>Automated enforcement and compliance.<br>Hierarchical decomposition of SLAs.                                                          |
| Data management                      | Data center data deduplication, petascale of structured and unstructured data; disks and tapes or backups; regulatory compliance.                         | Global deduplication, exascale largely unstructured data;<br>hierarchies of storage around NVRAM with disks at the<br>bottom; global compliance.                                                   |
| Integration of management components | Component integration at a single layer, local feedback loops; rapid deployment, configuration management and patching; orchestration of global services. | Choreographies and closed loops of loosely coupled<br>domains addressing power, performance, availability,<br>etc. individually and with tradeoffs (e.g. power-<br>performance for power capping). |
| Quantifying manageability            | Checklist of management functions,<br>documentation, time and steps to manage objects<br>and services.                                                    | Measuring Quality of Management (QoM), elasticity of management (matching manageability capabilities to those of functionality supported), ease of management.                                     |

Table 2. Summary of State of the Art and Research Direction of Cloud Management.

Table 3, Summary of Trends Impacting Future of Cloud Management.

| Layers of the Stack                       | State of the Art                                                                                                                                                                                                                                                                                                     | Research Direction                                                                                                                                                                                                                                                                                                                                         |
|-------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Cloud users                               | Traditional Internet users, increased mobile access<br>limited from developed and emerging areas. Some<br>mash-up ability of limited number of users. Some<br>ability to customize and personalize accounts.                                                                                                         | Dominantly mobile access, development countries<br>growth towards 8B users, especially through mobile.<br>Extensive mash-ups through user composed services.<br>Extensive personalization and customization.                                                                                                                                               |
| Cloud services developers                 | Small number for traditional and mobile services<br>Few releases annually, careful testing<br>some service location awareness<br>New services through development.                                                                                                                                                   | Through composition, integration, large % of developers;<br>continuous roll-out of new releases, agile development.<br>Full location awareness; integrate with local services<br>available ubiquitously.                                                                                                                                                   |
| Cloud management<br>operators             | Cloud and Cloud service operators (small %)<br>increasing updates to services mobile devices<br>some high level dashboard, analytics<br>reporting, some prediction.                                                                                                                                                  | Merging role with Cloud developers (large %)<br>frequent updates to mobile services, access devices,<br>detailed business dashboards, visual analytics<br>what if analysis, prediction business outcomes.                                                                                                                                                  |
| Hardware and <i>its impact on support</i> | Disks, early adoption of SDDs, 10Gb/s Ethernet,<br>early adoption of optical interconnect,<br>16-24 Core CPUs, 100,000 server data centers,<br>Air cooling, very limited use of water cooling,<br><i>high resource redundancy, reactive and delayed</i><br><i>support, field engineers, complex software repair.</i> | NVRAM adoption, broad optical interconnect,<br>deployment, 1000+ Core CPUs, with sophisticated,<br>photonics off-on chips, 10 <sup>12</sup> + server data centers,<br>ambient cooling (commodity), liquid cooling (high end),<br><i>self-healing, proactive support, customer self-repair,</i><br><i>repair moving up the stack, restartable services.</i> |