

Distributed Cloud Operating System Development at Academia Sinica

Eric Yen

Academia Sinica Grid Computing Centre (ASGC), Taiwan

E-mail: eric.yen@twgrid.org

Felix Lee

Academia Sinica Grid Computing Centre (ASGC), Taiwan

E-mail: felix@twgrid.org

Wei-Jen Chang

Academia Sinica Grid Computing Centre (ASGC), Taiwan

E-mail: weijen.chang@twgrid.org

Syue-Yi Liaw

Academia Sinica Grid Computing Centre (ASGC), Taiwan

E-mail: syueyi.liaw@twgrid.org

Shu-Ting Liao

Academia Sinica Grid Computing Centre (ASGC), Taiwan

E-mail: ting@twgrid.org

Chun-Wei Shen

Academia Sinica Grid Computing Centre (ASGC), Taiwan

E-mail: waynesan@twgrid.org

Jing-Ya You

Academia Sinica Grid Computing Centre (ASGC), Taiwan

E-mail: jingya.you@twgrid.org

Jhen-Wei Huang

Academia Sinica Grid Computing Centre (ASGC), Taiwan

E-mail: jwhuang@twgrid.org

Abstract

Building an *efficient* data center to meet the requirements of big data analysis is the goal of the Distributed Cloud Operating System (DiCOS). Our vision of DiCOS is to enable the support of scientific research by order-magnitude improvements in processing capacity, time-to-science, and data center efficiency. The first production DiCOS

(version 1.0) was released in 2013 and has been deployed at selected sites that support physics, earth science and other scientific fields, leveraging the Worldwide LHC (Large Hadron Collider) Computing Grid (WLCG) distributed computing infrastructure.

International Symposium on Grids and Clouds (ISGC) 2014
Academia Sinica, Taipei, Taiwan
23-28 March, 2014

POS (ISGC2014) 029

1. Introduction

The goal of big data analysis is to discover new knowledge and make greater advantage of large data in much more effective ways. Based on the experience from the Worldwide LHC (Large Hadron Collider) Computing Grid (WLCG), the largest academic big data project that has processed on the order of 100 petabyte data annually since 2011, it is impossible to achieve these (big data) goals without the right infrastructure and computing model. The internet-wide distributed infrastructure generalized from the practices of WLCG and many other e-Science applications has been able to reliably support more than 1.5 million jobs a day, together with over 100 petabytes per year data movement and access since 2013.

Efficient data analysis capability at the scale of hundreds of petabytes is the most demanding requirement in the early 21st century, not just the academic activities, but also for the advancement of industry and our daily lives. These big data analysis needs drive the requirements for smart data center(s), distributed computing infrastructure, innovative middleware, advanced networking as well as application integration. The vision of DiCOS is to accelerate scientific discoveries by leveraging order-magnitude improvements in data analysis, scale of computation, and resource federation efficiency. Thus allowing the users to focus on research and findings without resource limitations. DiCOS enhances the Distributed Infrastructure with Cloud services for multi-disciplinary scientific applications. In dealing with big data, a distributed infrastructure usually wins by the advantages of parallelism, throughput, safety, and availability etc. In fact, the flexible computing model also can avoid moving big files around and wasting time, network bandwidth, and other costs.

The strategy of DiCOS is to focus on the requirements of domestic academic communities first to accelerate their research computations, covering the domains of earth science, environmental changes, life science and physics primarily. Thus, the system technology and intelligence is continuously driven by user requirements and real applications from multidisciplinary user communities.

DiCOS is a generic big data analysis infrastructure overlay on the existing internet-wide highly scalable distributed infrastructure for multidisciplinary application requirements. Together with intelligent and flexible data processing capability, as well as the optimized system efficiency, DiCOS is designed to federate distributed scientific cloud resources and support customized workflow to strengthen the “User Pulls and Technology Pushes” model.

Examination of the daily usage patterns of various scientific applications, has allowed us to implement an effective and practicable system. Our chief approach is to identify and then develop the technologies to that are required to implement the required functionalities of the user communities. The distributed computing system is the most viable solution for us in Taiwan to make use of the global resources around the world for big data problems collaboratively, as the WLCG has successfully demonstrated. In

short, the development of DiCOS is to move forward the WLCG technologies for the local user community needs which are not well focused in Taiwan now.

A limited number of institutes in Taiwan have experience in operating and optimizing large-scale production distributed systems. The goal of DiCOS is to proactively extend both the user domains and the remote resource centres, based on the innovative energy-saving single rack as the building block. Quantitatively, the DiCOS is striving for a worldwide distributed multi-disciplinary high throughput computing system capable of running 10M jobs/day, mobilizing 1PB data over thousands of federated sites reliably.

2. Related Study and Our Experiences

Academia Sinica Grid Computing Centre (ASGC) is the only WLCG Tier-1 Center in Taiwan and has been operating since 2005. WLCG was the first and remains the largest big data analysis facility in the world with its Internet-wide scale, 1.5M jobs/day capability, 15 GB/s constant average data transmission rate mobilizing 180 Petabyte data among 300 resource centers in 40 countries over the past three years. By leveraging the successful experiences of WLCG and the collaboration with CERN, ASGC strengthens the multidisciplinary big data analysis competence of the new generation distributed computing infrastructure and innovating the core technologies to advance the scientific frontiers by working closely with various domain research groups.

Distributed computing infrastructure has proven to be an effective big data analysis platform by the WLCG. Moreover, the scalability, reliability (by fault tolerance and replication for example), and throughput are the outstanding advantages than other systems. During 2010 to 2013, ASGC processed over 33 petabytes of data in and out of Taiwan Tier-1 and Tier-2 centers. The inbound network traffic reached 16Gb/s in 2013 as shown in Figure 1. ASGC is now aiming at the research and development of the building block of a smart data center as well as the distributed computing infrastructure and technology. The hardware part is to design the fanless Single Rack Data Center (SRDC) with power usage effectiveness (PUE) < 1.2 . The SRDC is also a UPS-free facility with power supplies that can tolerate power outages of 30 seconds, and increases the power efficiency by up to a factor of two. Within 30 seconds of a power outage, the fast switch-on power generator for the data center takes over. Such improvements represent a significant savings in power consumption and the consequent CO₂ emission reduction. As an additional benefit, the usual acoustic noise level of over 90dB in a data center can be reduced to almost 0dB. The software part is to develop the Distributed Cloud Operating System. A generic big data platform supporting many big data applications running on the system at the same time could be built by taking advantages of dispersed resources over the Internet. Resources are easily federated and effectively provisioned in the form of on-demand services. Identify the performance bottleneck by the efficiency analysis from the data center power and thermal measurement to the SRDC and the core components of DiCOS are of highly essence to the ASGC.

From 2012, the European Grid Initiative (EGI) launched the EGI Federated Cloud project. The vision is to “Establish by 2020 a distributed open compute and data infrastructure comprising a 10M Core Federated Cloud and 10 Exabyte of Federated Cloud Storage across Europe that is able to support the data analysis activities of all researchers within the European Research Area” [1].

DiCOS v1.0 was deployed at ASGC and National Central University supporting large-scale AMS experiment analysis and engaged various fields research groups. About 1.4M jobs using around 1.4M CPU-core Hours of AMS experiment analysis were completed and 600TB data were transferred between Taiwan and CERN by DiCOS during 2012 and 2013. Except for the collaborative development of the core technologies with CERN, Industrial partners had also been engaged in accelerating the hardware and software development.

Resource Utilization at ASGC

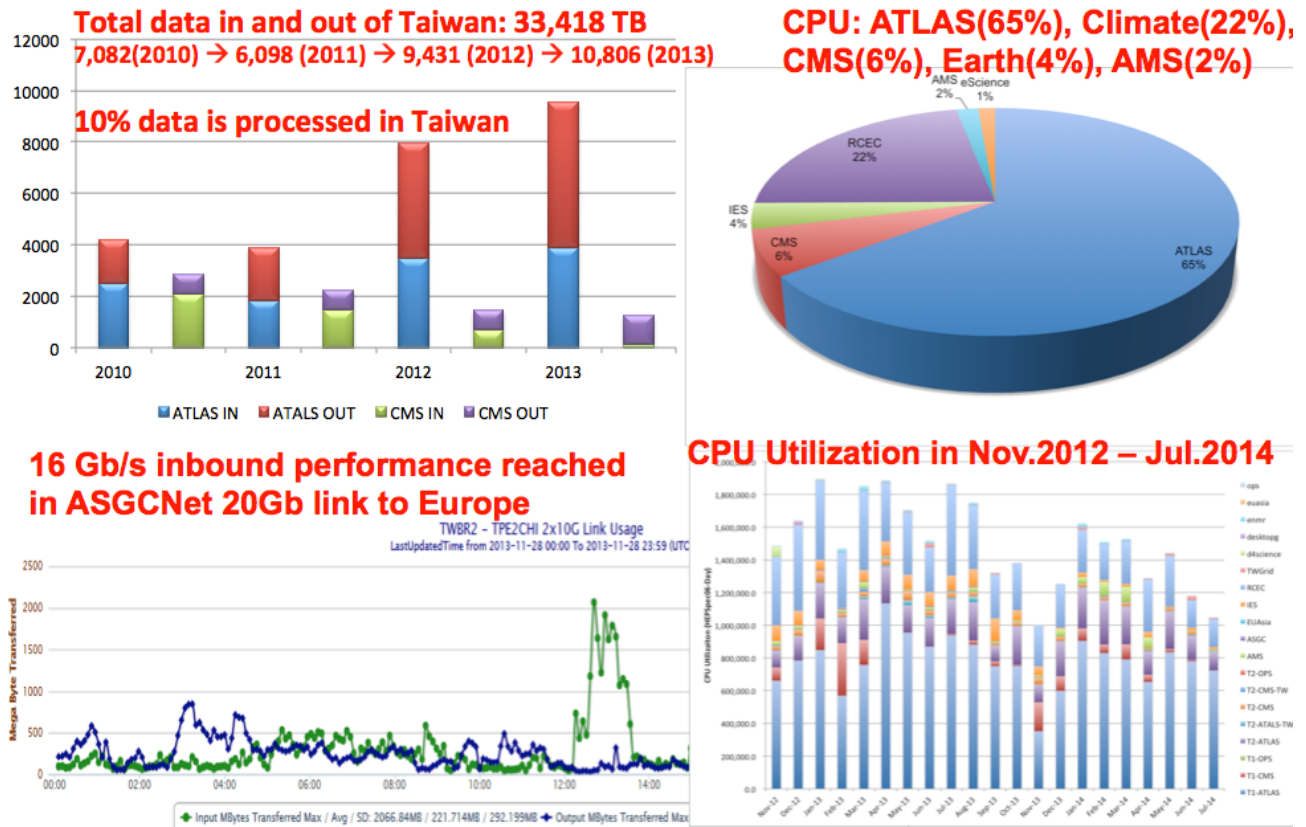


Figure 1. Based on experiences of development and operation of the worldwide distributed computing infrastructure in the past decade, ASGC is extending the production DCI for broader disciplines by DiCOS.

POS (ISGC2014) 029

3. Approach and Implementation

ASGC proactively evolves both the research domains and the efficiency of the applications by the experience from users and the advancement of core technologies. DiCOS provides an advanced distributed computing environment for the daily requirements of scientific big data applications in various domains. The result is a generic scientific cloud that allows many big data applications to be running on the system at the same time. Research and development, user community engagement, and international collaborations are the essential strategy to success.

To attain order-magnitude improvements, the first task was to design and build the new distributed infrastructure, based on the DiCOS infrastructure and the fanless single rack data center building block, at the collaboration sites in Taiwan. Collaboration sites are typically scientific user institutes or the resource providers that are able integrate local resources with ASGC through 10Gb Ethernet. Users are able to use any available resources in the distributed infrastructure without knowledge of the resource availability beforehand. Users have the highest priority on their local resources while their jobs are allowed to overflow automatically to the distributed infrastructure when there are no resources available at the users home institute. System metrics are recorded in order to allow the identification of performance bottlenecks so the system efficiency can be improved and operation costs can be reduced whenever an issue is identified and resolved.

DiCOS federates distributed resources of various scales consisting of core components such as the intelligent data management, distributed job management, pilot factory, Web user interface, Cloud framework and information system. Core software had been used by the WLCG connecting hundreds of sites and processing over 100 petabyte data every year. Architecture of DiCOS is illustrated in Figure 2.

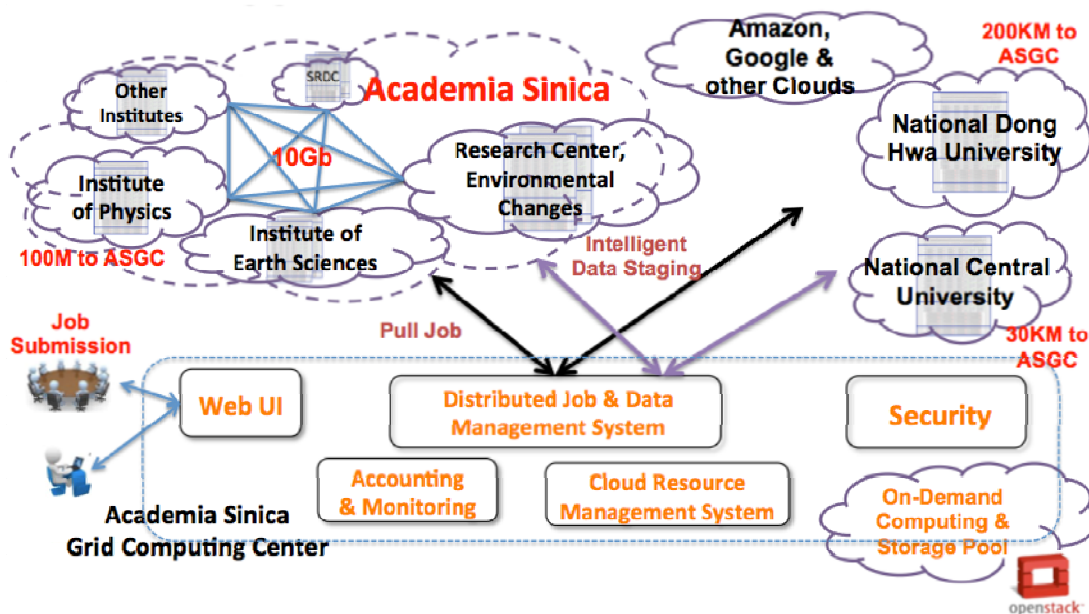


Figure 2. Distributed Cloud Operating System (DiCOS) Architecture with the core components identified.

Effective use of distributed resources for large volume of jobs and minimize the job latency is the first goal of the DiCOS Distributed Job Management component. Job is running at the best available site and data is moved transparently across sites whenever needed. Jobs are taken care by the pilot factory by pulling jobs from the available computing nodes. Computing nodes is dynamically provisioned with automatic deployment of required application environment. Application workflow is customized according to the computing model, apart from the generic web user interface. Data availability decides how the jobs would be executed and intelligent staging ensures the minimization of data transmission overhead. Integration of the system information system is essential for the smart wrapper (implemented by the pilot factory) to find the best resource for jobs and also re-submit failed jobs automatically. Until there is no jobs in the queue pilot then ceases to exit. The web-based job submission and job management interfaces of DiCOS are illustrated in Figure 3.



Figure 3. Web-based Distributed Job Management of DiCOS. The left is the job submission interface; the right is the job monitoring interface.

Distributed data management (DDM) is the key to big data analysis. Logically, the storage hierarchy is composed of disks, pools, sites and federation. In a heterogeneous distributed environment, DDM provides the collaboration among federated storages on data discovery, data transmission, data deletion, data consistency, and application computing model integration. Unified data catalog and global name space are the basis for the DCI to locate the required data effectively. Every resource center could have different storage system for the local data management. Through standardized interfaces, the local storage system could respond to inquiry and support data access and scheduled transmission by multiple protocols such as http, GridFTP, xrootd, etc. The DDM also provides policy-based data replication management to ensure data consistency and reduce the complexity systematically.



Figure 4. Directly drag-and-drop operation over the web browser is provided by DiCOS Distributed Data Management

Cloud technology is integrated into the DiCOS to realize fast provisioning of guaranteed execution environment at sites. Site level virtualization provides service consolidation and decoupling of jobs and physical resources to maximize the resource utilization while reducing the management cost. Certified virtual machine images and application environments are deployed by the CernVM File System [2, 3] based on the definition of virtual organizations and supported by the repository of virtual appliances. By the pilot job framework, the virtual machines are activated when there is any job coming to the queue.

On the other hand, consolidated storage is one of the core components to the DiCOS Cloud system right as the networking. OpenStack [4] is serving as the Cloud software of DiCOS and integrated with the Ceph [5] filesystem to provide the generic storage services to various requirements. In this regard, application systems could get storage services in types of object, block or file system from the generic and consistent cloud storage without setting up their own as usual. Besides the separation of namespace and underlying hardware, we could also replace the RAID controller and architecture by distributed storage and benefits from the removal of single point of failure.

Web user interface is devised to support integrated scientific workflow and to reduce the usage complexity at first. In current DiCOS version, the single sign-on mechanism and proxy management are also incorporated. All the job and data operations are web-enabled and intuitive drag-and-drop style data management is also supported as the Figure 4. Users could submit jobs, check status, manage data and visualize results by the web browser in any device.

System efficiency is of the greatest concern for a production system such as DiCOS. The efficiency not just covers the power, thermal and resource performance, but also includes all levels of system parameter tuning and operation efficiency to ensure the system performance and the fast bottleneck identification. Monitoring and configuration management are the two fundamentals of the system efficiency. A wide coverage monitoring framework has been built for both DiCOS and also the ASGC data center.

Identification of key metrics and moving towards automatic control are current focus based on daily analysis of the system status. Automatic load balance and early warning are the expected outcomes from the intelligent monitoring system in the future. A Puppet framework [6] has been deployed at ASGC for dynamic configuration and automatic system deployment. Better solutions for the issue of correlation between the software release and Puppet module release are being designed.

Besides the scientific disciplines mentioned above, applications of DiCOS will be extended beyond academic domains, such as the Cloud Internet of Things (IoT) services, supporting healthcare management and analysis, etc. Big data processing on DiCOS could be used in the mining of mobile users behavior patterns. Efficient collection and analysis data from IoT end devices as well as the social networks based on the requirements are all potential application scenarios as proposed by our university partners.

4. Result and Evaluation

Based on 10-year distributed system development experiences and the collaboration with heavy usage communities as aforementioned, ASGC is developing the green data center with fanless racks as well as the DiCOS, serving as the building block of distributed computing infrastructure and expediting data-intensive research. In addition to the ATLAS computing, the current production DiCOS successfully support AMS (Alpha Magnetic Spectrometer) experiment computing finished 1.4M jobs using 1.4M CPU-hours in Taiwan during Oct. 2012 to Sep. 2013. 75% jobs were running by cloud resources managed by DiCOS. 580 TB data in 700K files were transferred in and out of ASGC for AMS at the same time. This large scale real application also validated the efficiency of the first version DiCOS. Summary of AMS Computing support by DiCOS is depicted as Figure 5.

AMS Computing Supported by DiCOS of ASGC

- File Transfer Service via GRID
- SRM endpoints at CERN and TW
- Duration: Sep. 2012 – Aug. 2013
- CERN to TW
 - Volume: **423.654TB**
 - Number of File: **598,815**
- TW to CERN
 - Volume: **158.979TB**
 - Number of File: **91,097**

	Cloud	Cluster
#Jobs	810,944	596,388
CPU-Hours	820,000	576,113

- 400K+ Positrons of a common source found by AMS in 2013.
- 30% data analysis done by ASGC.
- ASGC serves as the one of the primary analysis center of AMS.

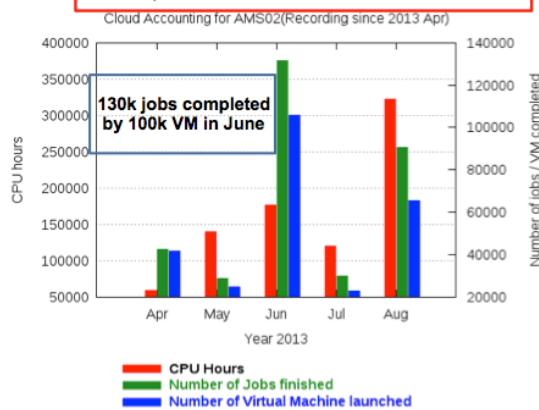


Figure 5. AMS computing support by DiCOS, around 30% AMS analysis jobs were done by DiCOS with both physical and cloud resources during 2012-2013.

Through the dynamic provisioning of VMs for AMS jobs, 100,000 VMs were generated by 130,000 jobs in DiCOS in June 2013. The additional overhead of AMS cloud based resources compared to physical machines, based on unit job execution time, is less than 5%.

In the first release of DiCOS, certificate-based single sign-on mechanism with flexible proxy life time management is equipped. Drag and drop Web user interface for distributed data management and job management are also provided. Dropbox-like data services is also supported.

User experience and requirements are the primary drivers for the improvements of DiCOS. ASGC keeps incorporating the right technology to fulfill the needs of various user groups according to the ICT evolution and the knowledge from partners. Apart from the AMS, DiCOS has been applied to difference fields such as the nuclear physics, polymer physics and bio-macromolecule physics, complex system, earth science and climate changes, and particle therapy etc. The DiCOS computing model has been extended to massive parallel computing framework and the license issue of commercial software is also investigated as summarized in Figure 6.

Applications Support in DiCOS

- Alpha Magnetic Spectrometer (AMS)
- Nuclear Physics (App: [GMC\(geant4\)](#))
- Polymer Physics & Bio-macromolecule physics
 - Computer simulations, Serial computing & batch submission
- Complex System
- Earth science and Climate changes
 - MPI & OpenMP (App: [gemb\(x\)](#))
- Particle Therapy: PTSIM
- Licensed (commercial) packages
 - License server at UI + group ACL

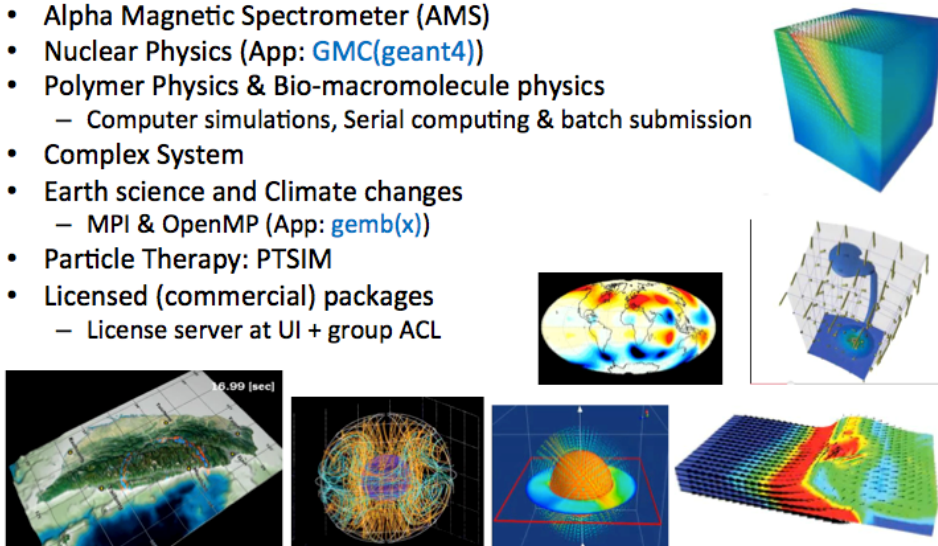


Figure 6. Widening the user communities and user domains are the key strategy of DiCOS development so that DiCOS could be evolving with better scientific computing performance.

5. Summary and Future Perspectives

Both the Distributed Cloud Operating System, the energy-saving single rack data center, the innovative cooling system, and the generic big data analysis platform are all novel in Taiwan. What is essential is that ASGC is integrating all these together and keep technology evolving by state-of-the-art advancement through international collaborations. With wider collaborations, we are not just extending the application domains but also outreaching to industrial partners. Fast and automatic deployment is one of top items of DiCOS in 2015. Using commercial resources as a DiCOS site then would be detailed explored.

The significant features of DiCOS are summarized below:

- A highly scalable and largest production-quality distributed computing system for generic big data applications in Taiwan has been established.
- It is capable of dealing with 100s PB data analysis in the throughput, limited only by the available network bandwidth.
- The building block (as of 2014) is a fanless and noiseless energy saving single rack with a PUE < 1.2.
- It is able to provide multiple disciplinary scientific big data analysis applications at the same time, eg, earth science, environmental changes, life science, physics, and astronomy.

Additionally, DiCOS paves a solid ground of Taiwan research infrastructure for the next decade, to support the systematic development of scalable and adaptive distributed computing infrastructure; to develop widely usable software tools for large scale and big data sciences; and to train HPC performance tuning and application development experts.

References

1. D. Wallom (2014), EGI FedCloud.
2. P Buncic et al. (2011), A practical approach to virtualization in HEP, The European Physical Journal Plus 126(1) 1–8 10.1140/epjp/i2011-11013-1.
3. J Blomer et al. (2011), Distributing LHC application software and conditions databases using the CernVM file system, J. Phys.: Conf. Ser. 331 042003.
4. T. Fifield et. al., OpenStack Operations Guide.
5. S. Weil, S. Brandt, E. Miller, D. Long (2006), Ceph: A Scalable High-Performance Distributed File System, OSDI 2006.
6. Puppet, <http://projects.puppetlabs.com>