
Cloud Infrastructure for Research Computing and Laboratory Environment

Bach Dániel

Budapest University of Technology and Economics E-mail:
bach.daniel@cloud.bme.hu

Geist Éva

Budapest University of Technology and Economics E-mail:
geist.eva@iit.bme.hu

Guba Sándor

Budapest University of Technology and Economics
guba.sandor@cloud.bme.hu

Abstract: With the cloud technologies the opportunity expanded to develop further the computer supported learning and education systems. There are several influencing factor of an education oriented cloud infrastructure. One of them is that the users, consumers are just regular computer users without deep knowledge of a cloud or IT systems. Virtual machines are widely used in training environments, but the technical details can be discouraging for the end users. Although, there could be several complicated scenario for the advanced users that need to be handled. Other special behaviour is that the lifetimes of virtual machines for the student labs usually short, with burst deployment rate at the beginning of a class.

Using and configuring virtual machines is not the privilege of IT technicians any more. Our solution provides an easy-to-use interface for anybody to create and share virtual appliances with others. CIRCLE hides these details and provides a simple role based user interface for administration, enabling users to focus on their actual task at hand instead of the infrastructure.

CIRCLE is a complete and open source cloud solution that can be deployed with minimal effort on a single computer as well as on a larger cluster.

The main contribution of this paper is to introduce the basic concepts of our solution called CIRCLE and reports about experience from the 2 years usage in education.

Keywords: TODO, Metadata; Scientific Data Management; Data Sharing; Data Integration; Computer Supported Collaborative Work.

Reference to this paper should be made as follows: Rodríguez Bolívar, M.P. and Senés García, B. (xxxx) 'The corporate environmental disclosures on the internet: the case of IBEX 35 Spanish companies', *International Journal of Metadata, Semantics and Ontologies*, Vol. x, No. x, pp.xxx-xxx.

Biographical notes: Manuel Pedro Rodríguez Bolívar received his PhD in Accounting at the University of Granada. He is a Lecturer at the Department of Accounting and Finance, University of Granada. His research interests include

issues related to conceptual frameworks of accounting, diffusion of financial information on Internet, Balanced Scorecard applications and environmental accounting. He is author of a great deal of research studies published at national and international journals, conference proceedings as well as book chapters, one of which has been edited by Kluwer Academic Publishers.

Belén Senés García received her PhD in Accounting at the University of Granada. She is a Lecturer at the Department of Accounting and Finance, University of Granada. Her research interests are related to cultural, institutional and historic accounting and in environmental accounting. She has published research papers at national and international journals, conference proceedings as well as chapters of books.

Both authors have published a book about environmental accounting edited by the Institute of Accounting and Auditing, Ministry of Economic Affairs, in Spain in October 2003.

1 Introduction

With the cloud technologies the opportunity expanded to develop further the computer supported learning and education systems. There are several influencing factor of an education oriented cloud infrastructure. One of them is that the users, consumers are just regular computer users without deep knowledge of a cloud or IT systems. Virtual machines are widely used in training environments, but the technical details can be discouraging for the end users. Although, there could be several complicated scenario for the advanced users that need to be handled. Other special behaviour is that the lifetimes of virtual machines for the student labs usually short, with burst deployment rate at the beginning of a class.

Using and configuring virtual machines is not the privilege of IT technicians any more. Our solution provides an easy-to-use interface for anybody to create and share virtual appliances with others. CIRCLE [1] hides these details and provides a simple role based user interface for administration, enabling users to focus on their actual task at hand instead of the infrastructure.

CIRCLE is a complete and open source cloud solution that can be deployed with minimal effort on a single computer as well as on a larger cluster.

The main contribution of this paper is to introduce the basic concepts of our solution called CIRCLE and reports about experience from the 2 years usage in education.

1.1 Summary - The CIRCLE Cloud Computing overview

Service specification as a cloud base education system may be offered by a range of commercial and non commercial service providers, universities or public bodies direct to Public Sector clients - educational bodies as service consumers or via other Service Providers.

The service reference framework is adapted from a number of industry de-facto standards, categorization and extended with specific requirements to create CIRCLE. The paper cites industry frameworks and de-facto standards in general terms to support the specific technology solutions used to provide the unique features, requirements of CIRCLE:

CIRCLE vision is based on providing greater choice, flexibility and efficiency for professors, teachers, lab instructors, lab leaders. The specification of service will allow them

to ensure appropriate alignment of their business requirements, - lab environment with the services offered. The services offered is listed through the templates (see below) available in the CIRCLE.

The services covered by the CIRCLE fall into 2 main categories:

Platform as a Service (PaaS) which is used to provide a platform for creating the necessary lab, development environment for the professors, teachers, lab instructors, lab leaders and students. Infrastructure as a Service (IaaS) for hosting the lab environments and to rent processing, storage, networks and other fundamental computing resources.

The services could be deployed in a number of ways, the CIRCE is deployed as a Hybrid Cloud. It can fit in community cloud definition because it is used by teachers, students for educational and research purpose. On the other hand it could fit into private cloud definition as well, because it is used only by a certain educational community. Both cases it is important that the users can access the services from the intranet as well from outside the educational environment, campus (see below Network). The service to provide the necessary security requirement integrated with the University central user - teachers, students, professors, etc. - database and use a special ACL structure. (see below ACL).

The Templates will provide access, easy to use and maintain interface for the teachers to create and maintain the necessary lab, development, test environment for the students, end users. It provides help to create hundreds of parallel environments in the same time. Through the templates is possible to modify, allocate, request more resources in individual bases.

Also through the QoS can be seen the need for additional resources and added additional resources by the administrator. Services provided by Service provider in the learning environment must provide performance standards and availability which not necessarily to be measured on a formal way.

The services which in the Service Catalog, implemented through the Templates will not include service tariff, instead the CIRCLE has got a built Lease feature (see below Lease).

The above introduction describes how the CIRCLE meet the five key characteristics of the cloud model: on-demand self service ubiquitous network access location independent resource pooling rapid elasticity pay for use

1.2 Architecture

IRCLE is based on open source free softwares. It uses the free KVM based virtualization through libvirt. The communications between the modules is handled by Celery with AMQP protocol. The portal is based on Django a python web framework.

2 User interface and authentication

With the implemented approach, it is easy to use and learn the interface and functions without any understanding of the lower layers and mechanism of the system. We designed an easy to use, reactive and simple web interface.

To keep the interface up-to-date the system does not store any information in the database that is queryable from another component. Only a small interval cache helps the interface to reduce the overhead of the communication between modules. As result the virtual machine states and the representing web page is well synchronised.

The GUI is easy to learn because the page structure remains the same for each user but the amount of options is changing based on the user's permission level. Moreover managing

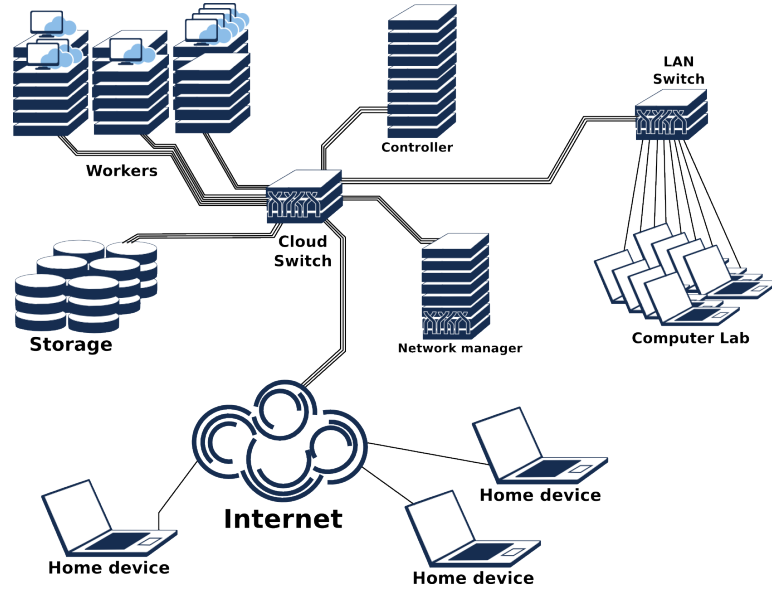


Figure 1 Architecture

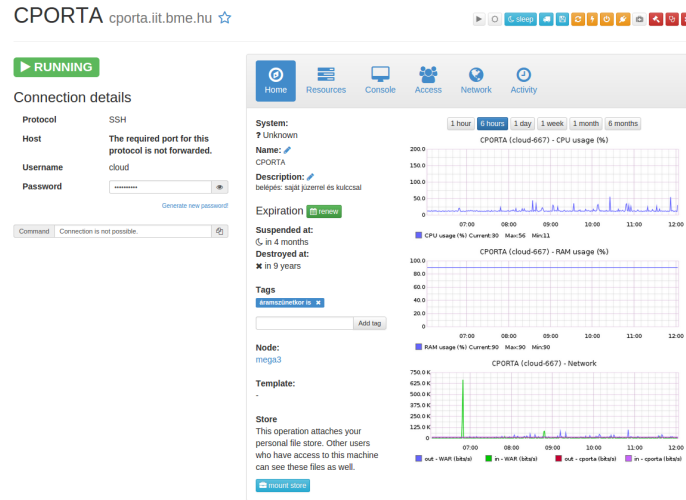


Figure 2 Dashboard view of a virtual machine

virtual machines is only possible in well defined scenarios. The user is not able to define inconsistent state of the configuration (i.e: a virtual disk as cd-rom). In our view the basic element is the virtual machine itself. For example the system does not allow users to define disk images without associating with virtual machines or templates. This helps to reduce the unused or forgotten disk images to save space on the storage.

The site supports both users and groups with permissions that makes easy to tune the access levels for large group of users as well. The administrators are even able to switch

accounts to view the page as the user. This really helps solving problems that only happens in the user's own context.

The SAML2 [2] based authentication module is connected to the campus wide used eduID [3] system which enables us to define default groups based on the eduid role (like students, teachers, etc) with the proper permissions before they are even logged in. Moreover the system creates groups based on the lectures that the user applied for. It promotes automatically the lecturers of the group to be owners as well.

To help the users manage their own resources CIRCLE uses notifications to communicate the changes in the system. Users can subscribe for a daily email that aggregates all the notifications in one message. For the rarely used functions - like template creation - we also included tutorials/wizards that guide through the whole process.

3 Templates and Requests

A common scenario in the education is when the instructor prepares a virtual machine for later usage in the class. CIRCLE has a simple templating system that allows the lecturer to customize the virtual machines and share it with users or groups. Creating template is possible from every virtual machine instance. After the necessary modifications (eg. installing a software module) on the virtual machine are made the instructor simply have to click to save it as a template. With the proper permissions teachers can tune the virtual machine resources as well. They can set the number of CPU cores, the memory size and they can add, delete or resize disk images. Advanced users can manage networks and even append native instructions for the virtualization API (libvirt). After saving the customized template the teacher is able to share it with the students, users and groups for self-provisioning.

Filter by status: [ALL](#) [PENDING](#) [ACCEPTED](#) [DECLINED](#)

ID	Status	Type	User
8	❌ DECLINED	🕒 Lease request	👤 Artúr Manó Marschal (SAGBE5)
7	⚠️ PENDING	📦 Resource request	👤 Balázs Ludmáry (ZG85W7)
6	❌ DECLINED	📦 Resource request	👤 Kamilla Tóth (SICFYZ)
5	✅ ACCEPTED	🔑 Template access request	👤 Kolos Koblász (F4MJSQ)
4	❌ DECLINED	🔑 Template access request	👤 Kolos Koblász (F4MJSQ)
3	❌ DECLINED	🔑 Template access request	👤 Vilmos Nagy (VIRNN9)
2	✅ ACCEPTED	🔑 Template access request	👤 Ferenc Varga (D89ZWT)
1	❌ DECLINED	🕒 Lease request	👤 Dániel Bach (J11M92)

Figure 3 List view or requests

Due to the limitation of the resource our students don't have permissions to start virtual machine until they granted by the teacher. However there are cases when the students need to ask for a template or more resource or even for more lease time for their project. CIRCLE integrates a requesting system that allows users to ask for more resource. The system administrators can easily accept or decline these requests on the dashboard and the requested feature automatically apply. This reduces the time for both asking for resource and the setting of proper parameters.

4 Sharing with ACL

To share a resource CIRCLE has a simple role-based ACL system. One can share everything including templates, virtual machines and networks with any users or groups. CIRCLE offers three default access level for every resources. The least privileged one is called User. 'User' is only able to use the shared resources. That means he can start a virtual machine from the template, add available virtual network or see the login credentials for a virtual machine. The next level is called Operator who is able to share the virtual resources for others as well. He can add and revoke User and Operator access rights. 'Operators' are able to do modifications on the appliances. They can restart, suspend, wake up or stop them. The highest level is 'Owner'. 'Owner' level provide the privileges like the original owner of the instance. 'Owners' can do every available operations on their own resources. They can share and destroy them as well. The only difference between 'Owner' level access and the original owner is that the original one can not be demoted. To remove an owner the ownership must be transferred to another user. With this small ACL system it is easy to share resources in a small group with the proper rights.

Permissions




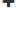
Who	What	
 Dániel Bach (JI1M92)	user	<input type="checkbox"/>
 Ádám Dudás (WG9114)	operator	<input type="checkbox"/>
 Sándor Guba (TFDAZ6)	owner	<input type="checkbox"/>
 <input type="text" value="Name of group or user"/>	user	

Figure 4 List view or requests

CIRCLE associates permission object with every functions. It is possible to create well defined groups for administering different parts of the system.

5 Leases - VMs lifecycle

Due to the unique characteristic of education the users are not forced to pay attention to minimize the usage of the system's resources. In business environments the users are charged for each virtual machine they started, even they use it or not. To prevent the system from using up all resources every instance has a predefined lease time. The lease defines two values: suspension and destruction time. For example a student laboratory takes 4 hours, and one of the lease value set up for 5 hours until suspension and 1 week for destruction. The lease can be extended by the user at any time via the renew function. In this way the student can continue his work at home, if he or she wants, but the unused VMs automatically deleted by the manager module after a certain time based on the set up.

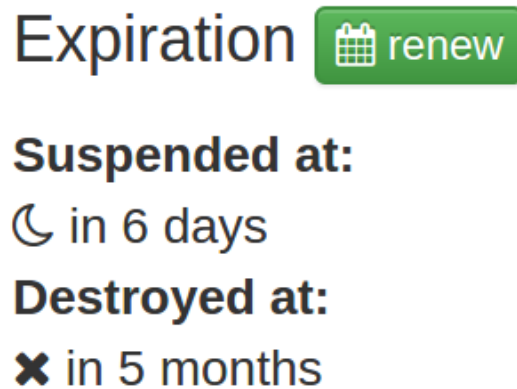


Figure 5 List view of requests

To follow the history of a virtual machine every operation triggers an activity event. The whole history of a virtual machine is tracked by activities. It stores the operation type, date and user. It also tracks the state modification created by the monitoring system.

Users with the proper permission can access to the virtual machine console as well.

6 Agent

Our cloud manager uses a small agent written in python on the guest machine. This agent helps the configuration and contextualization of the virtual instance. The agent copies the user's SSH keys, it changes password and changes the hostname. The agent does the network configuration as well. It provides feedback to the user when the computer's operating system is ready. The agent includes an alert function that notify the user inside the virtual machine when the VM is about to be suspended or destroyed. The agent is available on both Windows and Linux systems. It uses the virtio serial port to communicate between the host and the guest operating system therefore it works without network as well. The agent can update itself on the same channel as well.

7 Network

CIRCLE comes with a complex network administration system. This system manages both the physical and the virtual machines as well. The system provides VLAN tagged networks to separate network traffics from each other. Advanced users can define virtual networks in seconds. The CIRCLE network provides DHCP and DNS services if needed. It also includes a simple network-filter based firewall and IPv6 support. The network module offers several way to connect the virtual machine to the internet. It is possible to use public or network address translated (NAT) IP behind the firewall. Using NAT address will automatically generate a random port number to the public interface for the SSH or RDP protocol. This way connecting to virtual machine is simple from both LAN and WAN access. To prevent

Activity

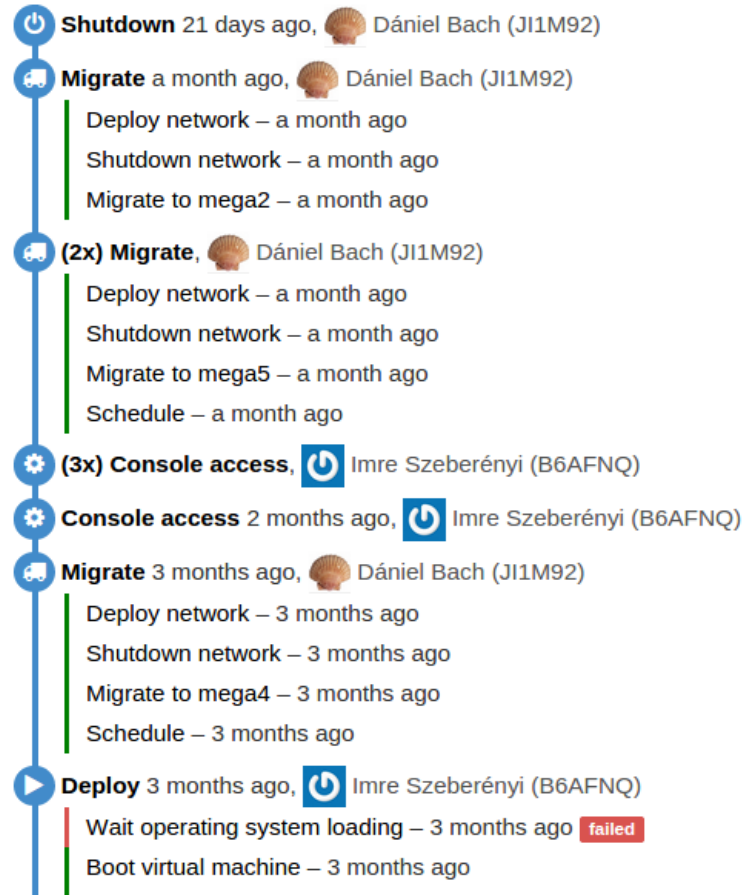


Figure 6 ACTivity history of a virtual machine

malicious or accidental abuse on the network there are several default ACL settings on the virtual network interfaces. Default it blocks the DHCP servers and the usage of other IP addresses than it's own.

Even users with the lowest permission level can perform some low priority modifications like port forward for the virtual machines. This is extremely useful for deployment with small public IP range but huge number of users (like a learning environment).

We are planning to release the network component as independent module. It helps the administering of small and medium physical deployments as well.

8 Storage

IRCLE is using shared storage to store the virtual machine images. Any shared filesystems like NFS can be used. There are two benefits using these kinds of storage. First, the virtual machines can be live migrated between the worker nodes and second the images can be

stored in the qemu copy-on-write file format. This makes possible of deploying virtual machines without cloning the base disk image (which can be huge). As a result it is possible to deploy lots of similar virtual machine at the same time in several seconds.

Saving a virtual machine will create a new image from the base and the differential image. This image can be used as base image for future deployment.

The CIRCLE storage manage the free space as well. You can define the free space in percentages. The images are not destroyed until the specified space is available. The garbage collector check it in a while and destroy the oldest images until there is enough free space again. Destroying a virtual machine on the dashboard will only move the image file in a trash folder. If the virtual machine's image is still on the storage (in the trash folder) it is possible to recover a deleted VM. Which is a useful feature if you forgot to renew some leases.

9 Easy installation

Installing complex systems with many modules are always troublesome. One of our main goal was to make the installation as easy as possible. We chose SaltStack a configuration management utility written in python. CIRCLE comes with SaltStack state files for each module. And with some “meta” installer for example an all-in-one package that installs a demo/developer machine. The only modifications need to be done in the SaltStack pillar files, that contains options, passwords and IP addresses. We are working on a small interactive script to make the installer a “next-next-finish” like software.

10 System admins

Managing the whole cloud can be done on the same user interface as well. There are several lightening of the management. It is possible to update the node via the web interface. In the background the salt installer can do all the needed modification (including git pull and managing the service restarts). Nodes have different states as well that help performing updates maintenance. Nodes can be Active, Passive, Disabled and Missing. The Active state means the fully functional node. The passive state is a functional node but the scheduler ignores it so no more virtual machine will start there. The disabled is a known missing or halted machine. And the missing means that the node is unreachable because of an error. Another useful operation is the flush. Practically it migrates all the virtual machine from that node to another one (recommended by the scheduler).

Updating nodes are easy as well. Based on our SaltStack installer the node can be updated from the Dashboard with 1 click. It will automatically checkout the new git repository and do the needed modifications as well.

11 Conclusion

The operated CIRCLE services point of view the basic requirement fulfilled and has a positive feedback from the users.

The basic Requirements of an education oriented cloud infrastructure: really easy to use, easy to maintain, no deep computer knowledge required Burst VM deployment, more

than 100 parallel is generated in seconds Wasting of cloud resources Main roles are teachers and students easy to use user interface, same interface is used by the maintainer TODO

From the perspective of users - teachers, students - are happy with the services and performance of the CIRCLE, fits in the education system very well because it helps the to be efficient, fast, more productive both for students and teachers.

12 Future work

We are working on a configuration management GUI integrated into CIRCLE ecosystem. The main concept is that there are predefined services with configurable options. The user can place services on a canvas and connect them through defined connection types. The user can save their workspaces and load if needed. When the architecture is finished it can be deployed on the cloud with one button. The software instances can be connected with “connections”. This way you can easily click together software architectures without deep knowledge of the software. The backing system would set up the software, complete the settings and care for minimal firewall as well. The manager’s OCCI interface is under development as well. Basic tests such as create and manage virtual instances are running successfully. The reporting portal for SLA and QoS will be further developed.

References

- Myneni, S. and Patel, V.L. (2010) ‘Organization of biomedical data for collaborative scientific research: a research information management system’, *International Journal of Information Management*, Vol. 30, No. 3, pp.256–264
- Bekelman, J.E., Li, Y. and Gross, C.P. (2003) ‘Scope and impact of financial conflicts of interest in biomedical research: a systematic review’, *JAMA*, Vol. 289, No. 19, pp.454–65
- Stein, L.D. 2008 ‘Towards a cyberinfrastructure for the biological sciences: progress, visions and challenges’, *Nature Reviews Genetics*, Vol. 9, No. 9, pp.678–688.
- NIH Statement on Sharing Scientific Research Data, http://grants.nih.gov/grants/policy/data_sharing/
- Policy for sharing of data obtained in NIH supported or conducted genome-wide association studies, <http://grants.nih.gov/grants/guide/notice-files/not-od-07-088.html>
- Network for Translational Research (NTR): Optical Imaging in Multimodality Platforms, <http://imaging.cancer.gov/programsandresources/specializedinitiatives/ntroi>
- Piwowar, H., Becich, M., Bilofsky, H. and Crowley, R. (2008) ‘PLoS medicine’, *Sept, No. 9, Towards a Data Sharing Culture: Recommendations for Leadership from Academic Health Centers*, Vol. 5.
- Birnholtz, J.P. and Bietz, M.J. (2003) ‘Data at work: supporting sharing in science and engineering’, *GROUP*, pp.339–348.
- Data Sharing & Intellectual Capital (DSIC) Workspace, https://cabig.nci.nih.gov/working_groups/DSIC_SLWG

- Getting Connected with caBIG, https://cabig.nci.nih.gov/getting_connected/
- Digital Imaging and Communications in Medicine (DICOM), <http://medical.nema.org/>
- The Cancer Genome Atlas (TCGA) Data Portal, <http://cancergenome.nih.gov/dataportal>
- National Biomedical Imaging Archive, <https://cabig.nci.nih.gov/tools/NCIA>
- caBIG: cancer Biomedical Informatics Grid, <http://caBIG.nci.nih.gov/>
- Biomedical Informatics Research Network, <http://www.nbirn.net/>
- Cyberinfrastructure for the Biological Sciences: Plant Science Cyberinfrastructure Collaborative (PSCIC), <http://www.nsf.gov/pubs/2006/nsf06594/nsf06594.htm>
- MIRC, <http://mirc.rsna.org>
- Foster, I. and Iamnitchi, A. (2003) 'On death, taxes, and the convergence of peer-to-peer and grid computing', *IPTPS'03*.
- Keidl, M., Kreutz, A., Kemper, A. and Kossmann, D. (2002) 'A publish & subscribe architecture for distributed metadata management', *ICDE*.
- Taylor, N.E. and Ives, Z.G. (2006) 'Reconciling while tolerating disagreement in collaborative data sharing', *SIGMOD*.
- Rader, E.J. and Wash, R. (2008) *CSCW*, pp.239-248, Influences on tag choices in del.icio.us.
- Halevy, A., Rajaraman, A. and Ordille, J. (2006) *VLDB, Data integration: the teenage years*, <http://portal.acm.org/citation.cfm?id=1182635.1164130>
- Doan, A. and Halevy, A.Y. (2005) Semantic Integration Research in the Database Community: A Brief Survey, *AI Magazine*, Vol. 26, No. 1, pp.83-94.
- Beynon-Davies, P., Bonde, L., McPhee, D. and Jones, C.B. (1997) 'A collaborative schema integration system', *Comput. Supported Coop. Work*, Vol. 6, No. 1, Norwell, MA, USA, pp.1-18, issn = 0925-9724.
- Wang, F. and Vergara-Niedermayr, C. (2008) 'Collaboratively Sharing Scientific Data', *CollaborateCom*, pp.805-823.
- Chin Jr., G. and Lansing, C.S. (2008) 'Capturing and Supporting Contexts for Scientific Data Sharing via the Biological Sciences Collaboratory', *CSCW*, ISBN 1-58113-810-5.
- W3C XML Query (XQuery), <http://www.w3.org/XML/Query>
- Wang, F., Hussels, P. and Liu, P. (2009) 'Securely and flexibly sharing a biomedical data management system', *SPIE*.
- Marlow, C., Naaman, M., Boyd, D. and Davis, M. (2006) 'Position Paper, Tagging, Taxonomy, Flickr, Article, ToRead', *Collaborative Web Tagging Workshop*.
- NCI Enterprise Vocabulary Services (EVS), http://ncicb.nci.nih.gov/NCICB/infrastructure/cacore_overview/vocabulary

NCI Enterprise Vocabulary Services (EVS), http://ncicb.nci.nih.gov/NCICB/infrastructure/cacore_overview/vocabulary

Bafoutsou, G. and Mentzas, G. (2002) 'Review and functional classification of collaborative systems', *International Journal of Information Management*, Vol. 22, No. 4, pp.281–305.

Bos, N., Zimmerman, A., Olson, J., Yew, J., Yerkie, J. and Dahl, E. et al. (2007) 'From shared databases to communities of practice: A taxonomy of collaboratories', *Journal of Computer-Mediated Communication*, Vol. 12, No. 2.

Myneni, S. and Patel, V.L. (2010) 'Organization of biomedical data for collaborative scientific research: A research information management system', *International Journal of Information Management*, Vol. 30, No. 3, June, pp.256–264.

myGrid, <http://www.mygrid.org.uk/>

Pike, W. and Gahegan, M. (2007) 'Beyond ontologies: Toward situated representations of scientific knowledge', *International Journal of Human-Computer Studies*, Vol. 65, No. 7, pp.674–688.