# GARR Cloud Storage  GARRBox

Cristiano Valli[a], Andrea Biancini[b], Fabio Farina[a], Fulvio Galeazzi[a],  Mario Reale[a],
Simon Vocella[a]

[a] Consortium GARR, Via die Tizii, 6,
00185  Roma, Italy
{cristiano.valli, mario.reale, fabio.farina, simon.vocella, fulvio galeazzi}@garr.it

[b] INFN Universita' di Milano Bicocca, Piazza della Scienza 3,
20126  Milano, Italy
andrea.biancini@mib.infn.it

**Abstract.** In this article an overview of Consortium GARR's overall strategy towards the provisioning of Cloud Services is given, with particular emphasis on the provisioning of a Cloud Storage service, GARRbox, currently in its prototypal phase. The basic architecture and the functional components of the system are described, and the future plans about it are highlighted.

**Keywords:** Identity Federations,  Cloud Storage,  Federated Access to Cloud Resources

## 1 Introduction

This paper describes the overall strategy towards the provisioning of Cloud service of the Consortium GARR (in the following GARR, for simplicity), the National Research and Education Network of Italy [1].
After providing the basic, guiding principles in the definition of a GARR strategy for Cloud, it describes in particular the prototype of a Cloud storage service known as GARRbox.  GARRbox is an ubiquitous access, multi protocol  Cloud Storage service prototype available to an limited community of GARR users for storing and sharing data files using both a Web-based interface, supporting Federated Single Sign on, and clients supporting the Amazon S3 protocol.

## 2 GARR overall strategy towards the Cloud

Looking back to the recent years, the preferred approach by users for accessing ICT resources has followed the policy of pursuing maximum efficiency and minimal required investment. The massive adoption of clusters and data centres technologies in the '90s and 2000, the collapse of the cost of connectivity, with always increasing available bandwidth, and, in addition, the fact that the average efficiency for using existing infrastructures is often very low, has brought to a revival of the computing model used until the '80s: the consolidation and sharing of resources used in centralized computing centres..

This principle, through the deployment of always improving virtualization technologies and Web 2.0 methodologies, has led to the emergence of the computing paradigm that today we identify by the term "Cloud". The lower cost of high-speed network connectivity with respect to hardware (with the electrical power and manpower costs necessary to manage it) made Cloud the winning approach compared to others.

From a user perspective, the Cloud has many attractive features:

• **On-demand Access:** Cloud ensures that user requests are promptly being answered.

• **Multi-modal, ubiquitous access**: users interact with the cloud using different protocols and devices accessing the data without knowing where they are physically stored or what the underlying provisioning technologies are.

• **Perception of unlimited resources**: Virtualization makes it possible to always ask for new resources.

• **Elasticity:** Resources are allocated and released automatically based on instantaneous workload, ensuring servicice continuity for users..

• **A well-defined business model:** Cloud resources are focused on the concept of utility provisioning; the economic model is then implicitly linked to the amount of consumed resources.

• **Little or no operating costs**: maintenance aspects of Cloud resources are delegated entirely to the Cloud provider, thus freeing users from costs and technical skills necessary to implement the same proprietary applications.

The Academic and Research worlds have perceived the potential of this paradigm. Different entities belonging to the GARR user community have advanced interest and inquiries about Cloud services, some trying to become service providers themselves.

It is therefore considered important for GARR to identify a course of action towards the Cloud, to better meet the demands of its community. This paper aims to describe GARR's approach towards the Cloud, highlighting models fitting both the mission and the available skills within the Consortium GARR and identifying the categories of Cloud service that could provide benefits to the whole community.

## 2.2 Key elements of GARR's strategy

As GARR is the Italian National Research and Education network, its main duty is the provisioning of a high capacity network backbone throughout Italy to interconnect the resource centres belonging to the Italian academic community to each other and to the rest of the world. Therefore, GARR is not aiming at becoming a provider of Cloud Computing resources, competing with existing public and private Cloud providers. However, GARR considers of strategic importance being able to add to its current services portfolio a higher level service devoted to the provisioning of Cloud Storage, for the immediate benefit of its user community, in particular, to start with, the e-Health community, already linked to GARR through several national (e.g.: National Health Ministry) and international (e.g.: the DECIDE project on e-Health) initiatives.

Moreover, its mission including enabling access to e-Infrastructures nationwide and at the international level, GARR intends to harmonize the current approaches its constituents organizations and research institutions are already adopting (or envisaging to adopt in future) , thus effectively ensuring the community of its users to be able to benefit from e-Science's advanced features, exploiting e-Infrastructures at all levels, even if different research bodies and institutions are providing and managing them. Another fundamental element GARR is taking into account is the effective exploitation of Identity Federations, providing users with single-sign-on to access a large variety of electronic resources.

In this sense, the key elements of GARR's strategy towards cloud can be summarized as follows:

- Acquiring direct, hands-on know how on the provisioning of Cloud Services
- Integration of Cloud provisioning to the national Identity Federation based on SAML2, called IDEM, also managed and provided by GARR
- Usage of standard, spread  protocols and platforms for accessing data, providing cloud resources

- Modular approach to the consuming and the provisioning of resource, i.e. being able to enable GARR's consitutent insitutions and universities to be able to easily become both providers of resources and consumers of federated Cloud services

The provisioning of a prototypal Cloud Storage service, called GARRbox, being able to respond to the basic requirements and needs of the (initial) community of e-Health.

More specifically, the requirement initially expressed by the reference e-Health users community are the following:

- Secure authentication mechanism
  - Whenever possible, based on the usage of Identity Federations
- Data integrity
- Ubiquitous access
- Dynamic Management of resources
  - Possibility to dynamically manage file access rights
  - Delegation of Management of subset of resources and Multi-tenancy
- Data encryption for data and metadata

# 3  The GARRbox Cloud Storage prototype

GARR started an internal pilot project on Cloud computing, normally referred to as GARRbox, in the third quarter of 2011, and the project is currently still ongoing.
The main goal of GARRbox is to acquire hands-on know how on Cloud Storage provisioning in view of setting up a service for the Biomedical e-Health Community.

In particular, experience has been gained through the pilot on:

- AAI based on the national identity federation (IDEM [2])
- Performances and main features of distributed file systems
- Providing different interfaces to users for accessing their data
- Business continuity / Disaster recovery

## 3.1 GARRbox Architecture

The prototype is based on a 3-tiers layered architecture: the lower layer is represented by the physical resources; the middle layer is the cloud filesystem and aggregation

layer, and the upper part is represented by the presentation layer and its user front-end. The current architecture of GARRbox is shown in Figure 1.
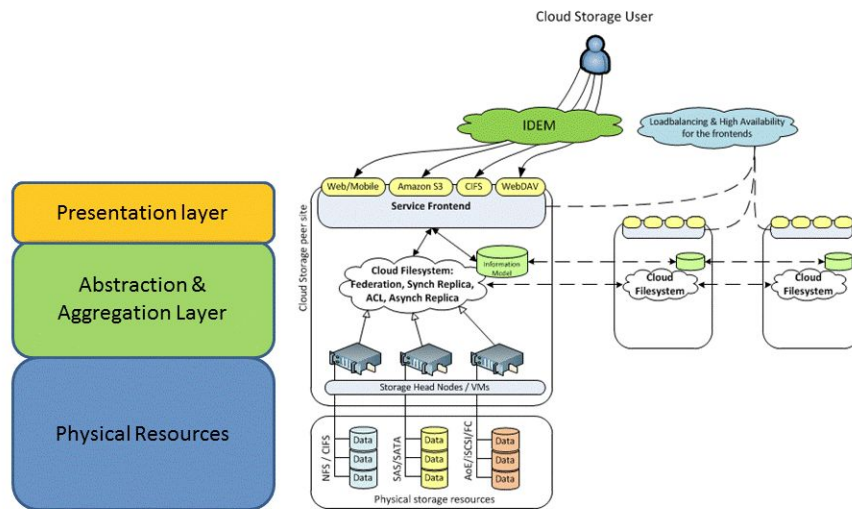


**Fig. 1.** The layered architecture of the current GARRbox protoype.

The prototype service is based on the usage of existing architectural components and tools, which have been integrated, customized and consolidated to provide a consistent and reliable service for the GARR user community. In particualar, the system is exposing both a web interface, a portal based on the Ajaxplorer tool [3], relying on the Single Sign on procedure provided via Shibboleth by the GARR Headquarters IDEM Identity Provider, , and an Amazon S3[4] endpoint interface which is originally based on CUMULUS, by the NIMBUS [5] project platform. Moreover, the aggregation layer is based on GlusterFS [6], the open source cloud filesystem .

The currently available interfaces are the Ajaxplorer based web/mobile interface (under shibboleth Single Sign on) and the Amazon S3 protocol interface, for which users do require secret and access keys, acquired during their registration phase to the GARRbox service. At the present stage, therefore, two independent AAI models are available: SAML2/Shibboleth-based for the Federated Access and Secret/Access key pairs for the AMAZON S3interface.

Essential features of the current adopted architecture are the following:

- Layered architecture
- Each layer self-contained within its interface boundaries
- Protocols and underlying technology within a layer replaceable without impact on others
- Resilience / High Availability

**Table 1.** Functional components / packages for the GARRbox prototype

| Functional Component | Tool | Reference |
|---|---|---|
| Web Front End | Ajaxplorer | http://ajaxplorer.info/ |
| Amazon S3 client | DragonDisk | http://www.dragondisk.com/ |
| Amazon S3 server | Cumulus/Nimbus | http://www.nimbusproject.org/ |
| Identity Provider/SP | Shibboleth | SAML 2 http://shibboleth.net/ |
| Information System | MySql | http://www.mysql.com/ |
| Cloud Filesystem | GlusterFS | http://www.gluster.org/ |
| Identity Federation | GARR IDEM | http://idem.garr.it |

**Physical Resources**

At the lowest level, the physical resources, the service currently makes use of 2 clusters of machines, both phisical and virtual ones, hosted at the GARR premises in Rome and the GARR Milano Bicocca site. Both clusters are provided with

- A machine acting as web front end/head node, where the Web access portal is installed (*garrbox.dir.garr.it* and *garrbox.mib.infn.it* respectively).

- Nodes belong to the GlusterFS pool, where GlusterFS bricks have been defined and assembled through GlusterFS into Volumes, providing the basic storage for users' data repositories
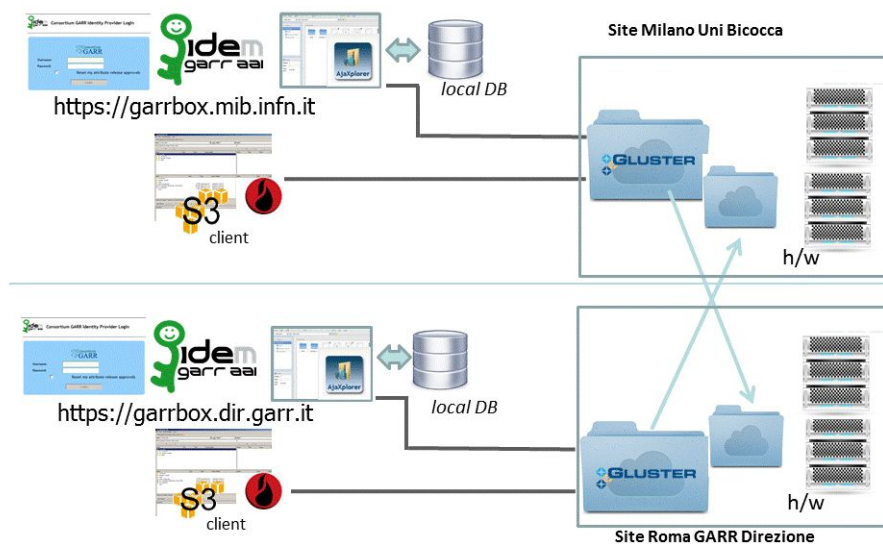
**Fig. 2.** The two-sites based provisioning of the GARRbox cloud storage service

**Aggregation Layer**

The aggregation layer is provided by GlusterFS deployed at two sites, GARR Headquarters in Rome and University of Milano Bicocca.

GlusterFS has been deployed and configured creating *distributed and replicated* type gluster volumes in both the Rome and Milano sites. Each cluster exports via GlusterFS to the corresponding front-end node 2 GlusterFS volumes : one acting as the main reposistory for users accessing the corresponding GARRbox front-end, and a second one, currently used for GlusterFS cross-replication between the Rome and Milan sites.

The system is therefore exploiting GlusterFS geo-replication of *distributed and replicated* Volumes between the two sites.

A detailed layout of the GARRbox layered architecture is shown in Fig.3.
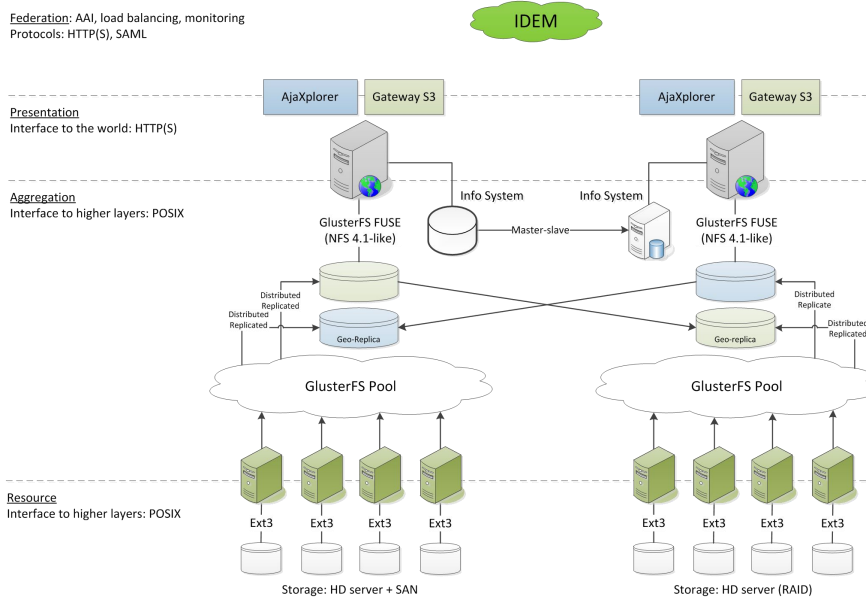
**Fig. 3.** Expanded layout of the GARRbox layered architecture

## 3.2 User Interfaces and provided functionality

As mentioned, GARRbox currently provides a web/mobile interface based on Ajaxplorer (for which a licenced mobile client exist), and an AMAZON S3 gateway endpoint. The system is currently available internally for GARR staff (around 20-25) users, each one with a quota of 10 GB storage space.

Users can make use of the available functionality provided by the system in the following way:

- Upload / download, search and manipulate data from client browsers and S3
- Sharing files and folders
- Single Sign On provided by IDEM IdP/SP
- Sharing files with anyone, shared folders by IDEM Attribute Principal Name (EPPN)

Furthermore, users can access the following functionality:

- Web authentication and access rights, access through credentials similar to the ones used by Amazon S3 clients (access key and secret key)
- Different quotas based on user, group and organization
- Local and geographical replication of files

- Soft (early warning) and Hard (further access denied) quota control mechanisms

The Dragon Disk S3 client is shown in Figure 4. It is one of the S3 tools supported by GARRbox and allows users to store and share files through access and secret keys, shipped to the user at the first login.
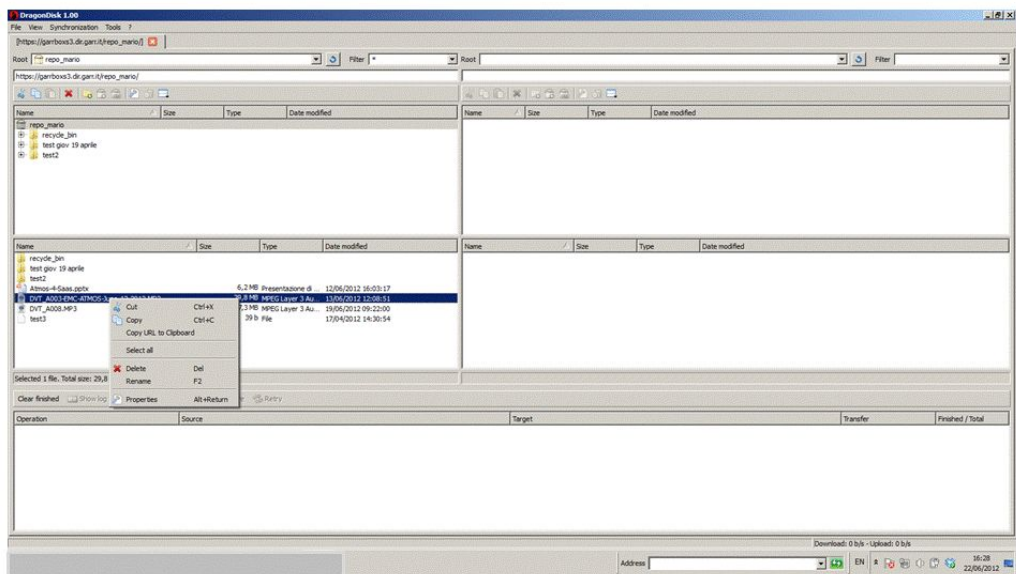


**Fig. 4.** The Dragon Disk based Amazon S3 client user interface

The Web based interface used by GARRbox is based on AjaXplorer. The AjaXplorer interface of GARRbox is shown in Figure 5. It currently allows users to store files on their personal repository, which is autocreated by the system at the first login. It also allows users to share files and folders, setting an expiration time for the sharing and generating a public links to be distributed to collaborators.
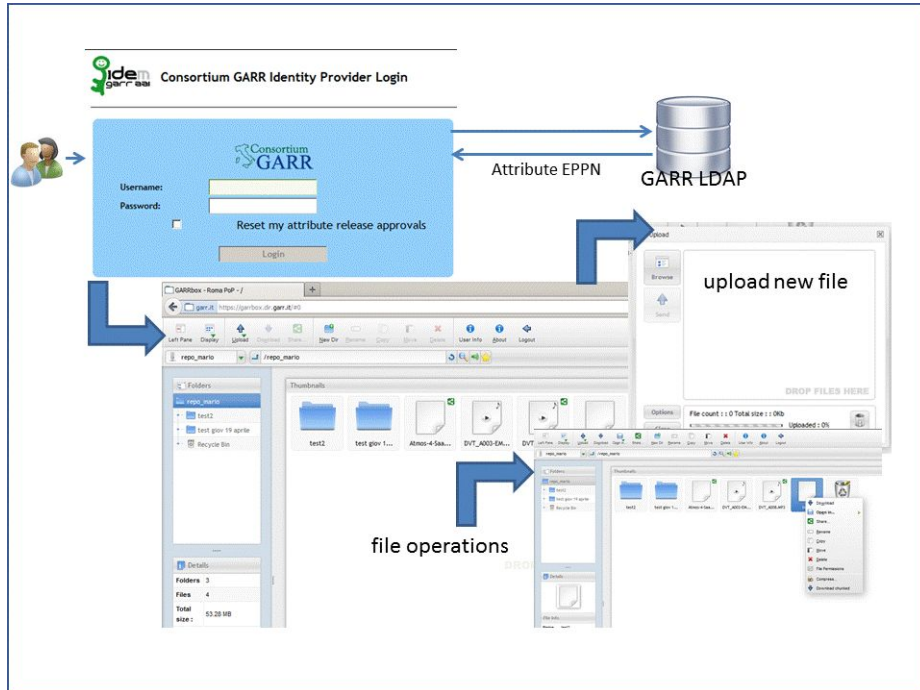
**Fig. 5.** The Ajaxplorer based interface to the GARRbox Cloud Storage service and associated user workflow, involving IDEM Single Sign On procedure.

## 3 Conclusions and Future Outlook

So far, the GARRbox system has performed reliably, proving to be able to respond to the first basic requirements of GARR users while using a Cloud Storage service. First tests of efficiency and resilience have yielded promising results.

The system is currently in its beta testing phase, after consolidation developments have been carried out.

Future developments will span different functional and architectural domains, and can be summarized as follows:

- **Identity Management**
  - Harmonization of credentials management: IDEM vs S3
  - Implement Authorization entirely based on IDEM / Shibboleth
    - add support for Administrative Roles
    - evaluate possible need for new Attributes and/or New Attribute Values
  - Integration of WAYF service [9]
- **Resilience in case of fault**
  - Transparent migration of user across backend sites
    - Switch over
    - Failover/failback
  - Information System DB implemented in master-slave / multi-master
  - High Availability for the front ends
- **Load Balancing**
  - Geographical DNS between the 2 front ends
  - Front end selection for users based on their IDEM attributes

Furthermore, further developments will involve the following required consolidation points

- Global **security assessment**
- **Presentation Layer**
  - Improving Web Interface
  - WebDAV interface
- **Aggregation Layer**
  - Client & Server side encryption
  - Management of complex metadata
  - Decoupling data from metadata
  - Multi-tenancy / Virtual resources management delegation
  - Assessment of alternative technical solutions:
    - NoSQL
    - Newer Gluster version or alternative distributed filesystems
- **Extend validation phase by users**
  - Extend beta users community

Future activities will include a deeper assessment of what done at the international NREN level and TERENA, beyond what done at project start, and the evaluation of synergies and collaborations worldwide.

Other long term developments will consider an improvement of the architecture to simplify the management and the adoption of the service. Furthermore, synchronization clients for Linux, Windows, OSX and device platforms will be developed. The integration of additional file access protocols like NFS and CIFS, and developments involving SSO bridging technologies spanning different (L2-L3) layers in the stack, will be taken into account.

## Acknowledgments

## References

1. GARR, the Italian Academic and Research Network  http://www.garr.it
2. The IDEM Identity Federation   https://www.idem.garr.it/
3. Ajaxplorer Web Filesystem browser http://www.ajaxplorer.info
4. Amazon S3  http://aws.amazon.com/s3/
5. Cumulus/Nimbus project http://www.nimbusproject.org/
6. GlusterFS cloud filesystem  http://www.gluster.org
7. Dragon Disk Amazon S3 client http://www.dragondisk.com/
8. Shibboleth  - SAML2 SSO implementation  http://shibboleth.net
9. Shibboleth WAYF  Service  http://www.switch.ch/aai/support/tools/wayf.html