# eGee

Enabling Grids for E-sciencE

# Deployment and Management of Grid Services for a data intensive complex infrastructure

*Álvaro Fernández  ( IFIC – Valencia )*

*5th User Forum*

*Uppsala, April 2010*

**www.eu-egee.org**

e-infrastructure

# Contents

- **Motivation for the talk**

- **Proposed Solution**

- **Computing and Storage Hardware**

- **Configuration and Management**

- **Developments**

- **Conclusions and Future**

**Enabling Grids for E-sciencE**

- **Several user communities require access to computational and storage resources in efficient manner.**
  - National Grid Initiative applications
  - Grid-CSIC
  - Atlas Tier2 and Tier3
  - Local users

- **Different kinds of applications and data access patterns**
  - Parallel and sequential Applications
  - Public and private data. Need for sharing policies
  - File size varies, sequential and random access

**Enabling Grids for E-sciencE**

- **To use grid technologies that we have been working with, and Hardware and Software configurations that alow to grow in the future.**

- **Storage and Network Hardware based on a NAS**

- **Lustre as a parallel backend filesystem**

- **Storm as a Storage Resource Manager  for our disk based storage**

- **Grid Services based on glite**

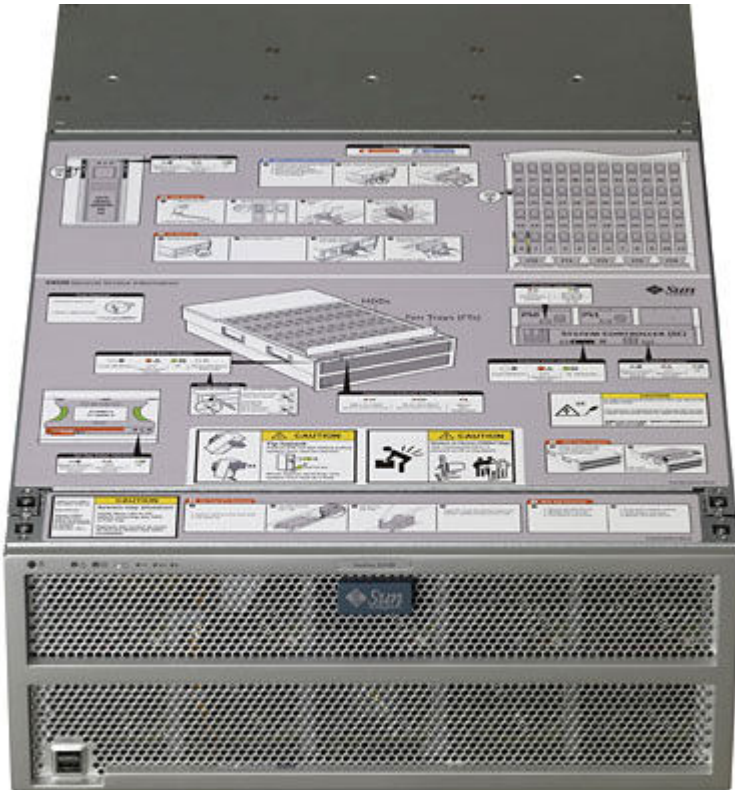- **Development of plugin, Tools and High level applications for user commodity**

- **Grid-CSIC resources**
- **Sequential:** 106 nodes
  - **2xQuad Core Xeon E5420 @2.50GHz**
  - **16 Gbytes Ram**
  - **2xHD SAS 134 BG RAID 0**
- **Parallel:** 48 nodes
  - **Same as before with Infiniband Mellanox Technologies MT25418 DDR 4x**

- **Total cores : 1704**

- **Atlas VO resources**
- **32 + 19 nodes HP DL160**
  - **2xQuad Core Xeon E5472 @3.0GHz**
  - **16 GB RAM**
  - **2xHD SAS 134 BG RAID 0**

- **Tape library for long time storage (with legacy castor sw)**
- **Disk servers to store online data (as part of Tier-2 requirements)**
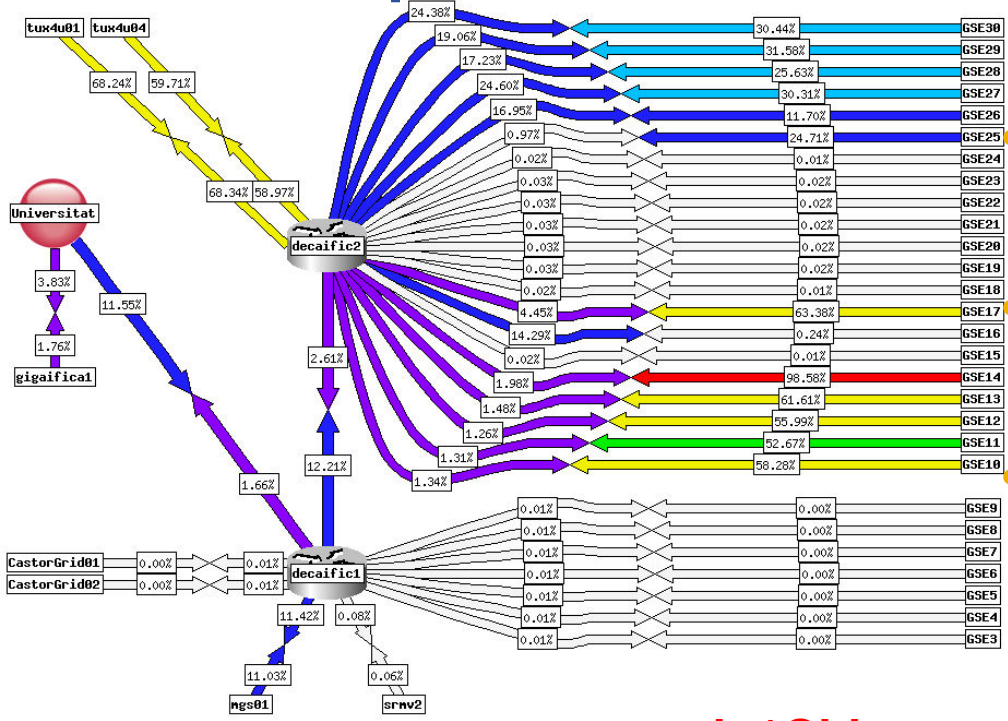  - Fast access and always available

- **Sun X4500/X4540 server**
  - **48 disks (500Gb/1Tb), with 2 disks for OS**
  - **Redundant power supply**
  - **4x1Gb Ethernet ports**
  - **Scientific Linux IFIC v4.4 with Lustre kernel**
  - **Ext3 system (volumes until 18/36 Tib)**

- **IFIC tested  best configuration:**
  - Disks in raid 5 (5x8 + 1x6). Usage ratio of 80%
  - Best performance with 1 raid per disk controller
  - Bonnie++ aggregated test results:
    - Write: 444,585 KB/s
    - Read: 1,772,488 KB/s
- **Setup:**
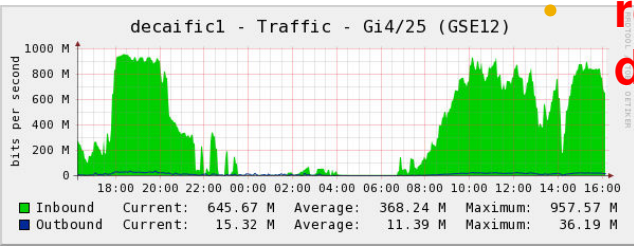  - We have 21 servers providing a raw capacity of around 646 TB

**Enabling Grids for E-sciencE**

- **Data Network based on gigabyte ethernet.**

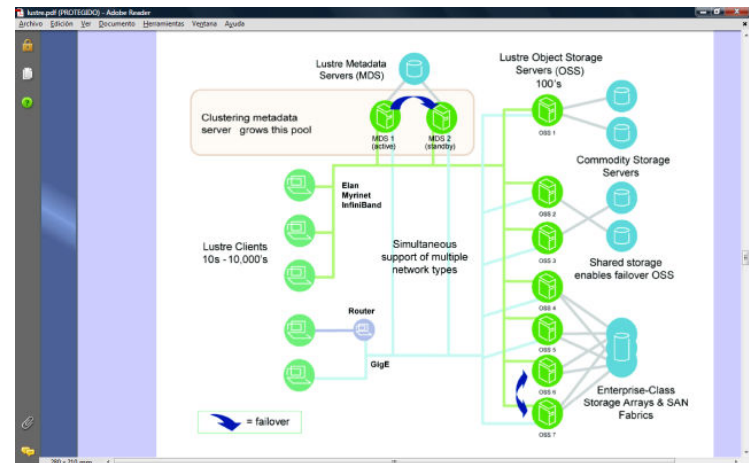- **10GB uplink to the backbone network**



- **Cisco 4500 – core centre infrastructure.**

- **Cisco 6500 – scientific computing infrastructure**

- **Data servers with 1GB connection. Channel bonding tests were made. We aggregate posibly 2 channels in future.**

- reach 1Gbit per data server

- **WNs and GridFTP servers with 1GB**
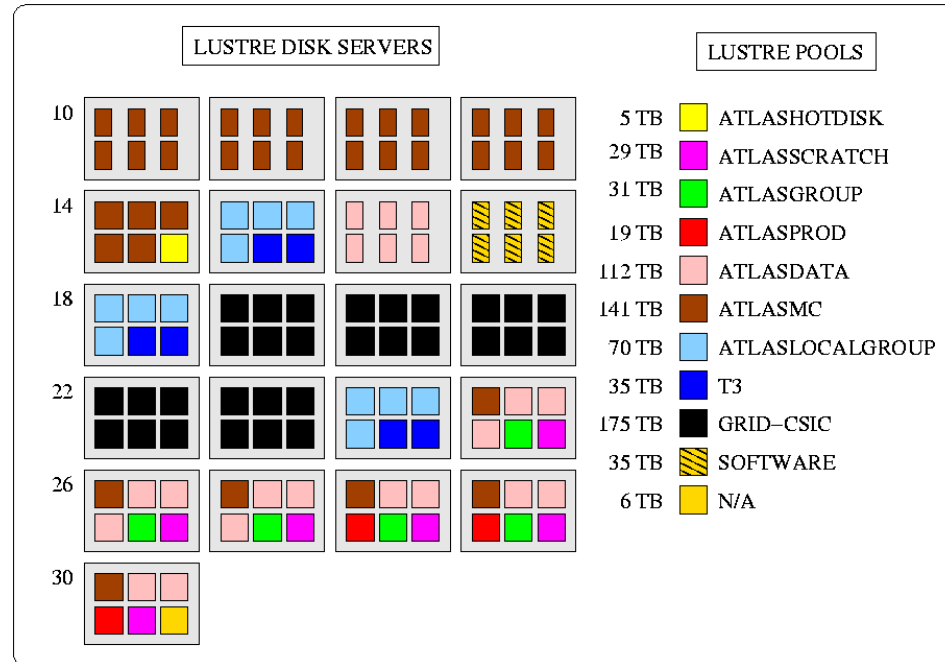
**Enabling Grids for E-sciencE**

- **We chose Lustre** because is a **Petabyte** scale **posix** filesystem, with **ACL** and **quotas** support. **Opensource** and supported by SUN/Oracle. In our installation with a single namespace ( /lustre/ific.uv.es):

- **MGS: Management Server** ( 1 per cluster ) + **MDS: Metadata Server**. ( 1 per filesystem )

  – 2xQuadCore Intel(R) Xeon(R) CPU
  – 32 GB.
  – Planned to add HA
  – Currently 4 Filesystems (MDT)



- **OSS: Object Storage Server**. ( Several per cluster guaranteeing scalability) 1 in every SUN X4500.

  – 6 OSTs per server (raid-5 5x8 + 1x6)

- **Clients**. ( SL5 kernel with lustre patches), mainly User Interfaces (UI) and Worker Nodes (WN) with configurations in RW and RO mode

- **With the Lustre release 1.8.1, we have added pool capabilities to the installation.**

- **Pools Allow us to partition the HW inside a given Filesystem**
  - Better data management
  - Assign determined OSTs to a application/group of users
  - Can separate heterogeneous disks in the future

- **4 Filesystems with various pools:**
  - /lustre/ific.uv.es Read Only on WNs and UI. RW on GridFTP + SRM
  - /lustre/ific.uv.es/sw. Software: ReadWrite on WNs, UI
  - /lustre/ific.uv.es/grid/atlas/t3 Space for T3 users: ReadWrite on WNs and UI
  - xxx.ific.uv.es@tcp:/homefs on /rhome type lustre. Shared Home for users and mpi applications: ReadWrite on WNs and UI
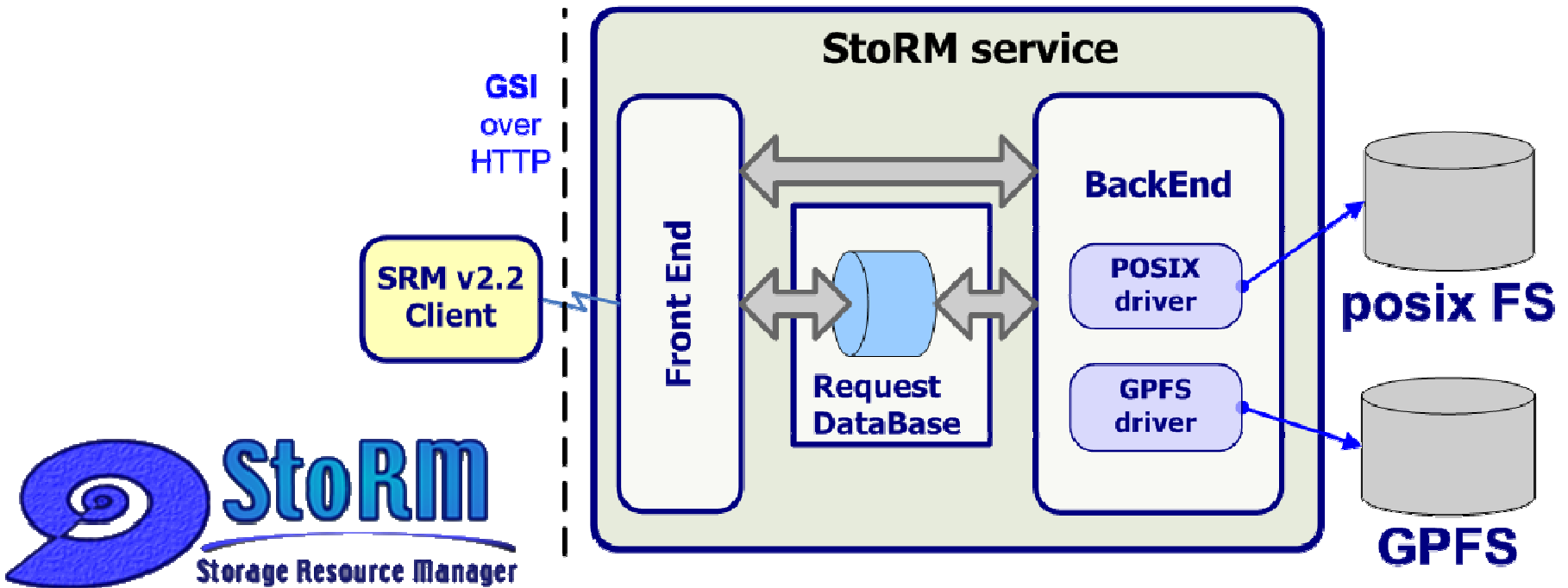
- **Quotas**
  - Applyed per Filesystem
  - For the Users in the Tier3 space
  - For the Shared Home Directory
  - Not being applied in general VO pools/spaces because all data is common for the VO
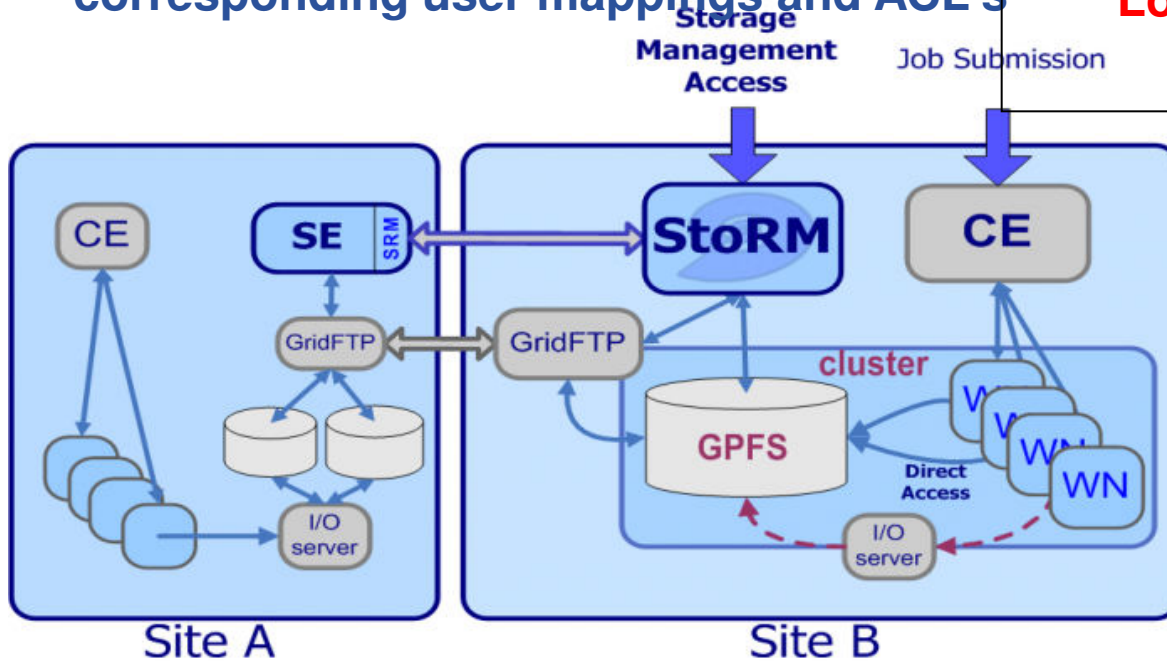- **ACLs**
  - Applied in the diferent VO pools/spaces to allow policies to be implemented (ie: only some managers can write, all users can read)
  - In the User and shared home directories
  - In the SW directories
  - This to be respected by the higher Grid Middleware Layers (STORM part of the presentations)

**egee**

- **Interface of the storage network with the GRID.**
- **For disk based Storage Elements**
- **Implementing SRMv2.2 specification.**
- **Allows accessing Lustre filesystem with the posix interface.**
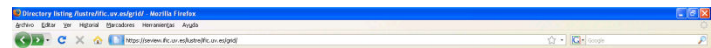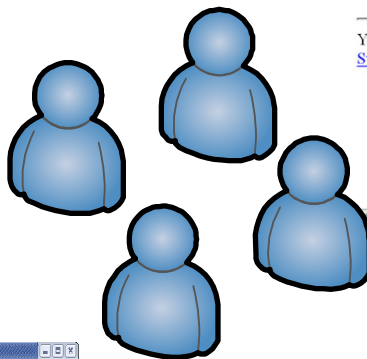
**Enabling Grids for E-sciencE**

- **Access like a local file system, so it can create and control all the data available in disk with a SRM interface. (BE + FE Machine)**

- **Coordinate data transfers, real data streams are transferred with a gridFTPserver in another physical machine. (Pool of gridfts)**

- **Enforce authorization policies defined by the site and the VO. (users are mapped to common vo account)**

- **We Developed Authorization plugin to respect local file system with the corresponding user mappings and ACL's**

**LocalAuthorizationSource available > Storm 1.4**



- # file: atlasproddisk
- # owner: storm
- # group: storm
- user::rwx
- group::r-x
- group:atlas:r-x
- group:atlp:rwx
- mask::rwx other::---
- default:user::rwx
- default:group::r-x
- default:group:atlas:r-x
- default:mask::rwx default:other::---

**Enabling Grids for E-sciencE**

- • **Besides posix CLI:**
- • **IFIC users additionally can access data through a WEB interface.**
- • **Ubiquitous and X.509 secured easy read-only access.**
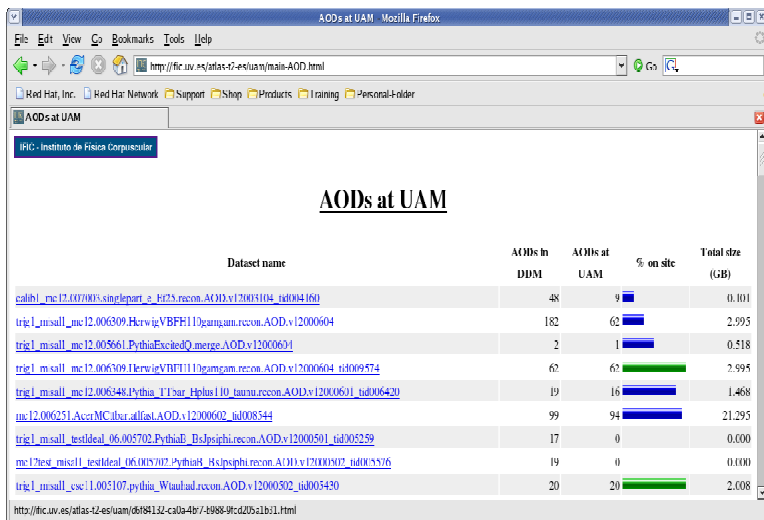- • **Developed with GridSite**



## Directory listing /lustre/ific.uv.es/grid/

[Parent directory]

| | | | |
|---|---|---|---|
| atlas/ | 4096 | 12:53 | 26 Jun 08 |
| dteam/ | 4096 | 00:23 | 28 Nov 07 |
| ific/ | 4096 | 17:15 | 19 Mar 08 |
| ops/ | 4096 | 10:49 | 23 Jul 08 |
| swetest/ | 4096 | 00:24 | 28 Nov 07 |

You are /DC=es/DC=irisgrid/O=ific/CN=Alvaro-Fernandez
Switch to HTTP . Built with GridSite 1.5.2

**WEB Interface to check datasets (Atlas groups of files) at the Distributed Tier-2**

**Enabling Grids for E-sciencE**

- **Presented a data management approach for multi-VO in a computer centre based on grid technologies and mature software.**
  - Tested well in stress tests (Step09,UAT)
  - more challenges to come, prepared to manage problems
- **Lustre gives good performance and scalability to grown next year until Petabyte scale.**
- **Storm presents our disk-based storage to the grid.**
  - Performed very well, and suits our necessities.
  - We would like it to continue not to being a complex system
- **In the future, deploy an implementation for HSM and Lustre:**
  - Follow CEA/Oracle work on this, these days: http://wiki.lustre.org/index.php/Lustre_User_Group_2010)

# Final words

**Enabling Grids for E-sciencE**

## More information:

[https://twiki.ific.uv.es/twiki/bin/view/Atlas/LustreStoRM](https://twiki.ific.uv.es/twiki/bin/view/Atlas/LustreStoRM)

[https://twiki.ific.uv.es/twiki/bin/view/ECiencia/WebHome](https://twiki.ific.uv.es/twiki/bin/view/ECiencia/WebHome)

[Alvaro.Fernandez@ific.uv.es](mailto:Alvaro.Fernandez@ific.uv.es)

## Thanks for your attention !