



# LCG Monitoring and Accounting

Dave Kant  
CCLRC e-Science Centre, UK

LCG Workshop  
Nov 2<sup>nd</sup>-5<sup>th</sup> 2004

Dave Kant  
D.Kant@rl.ac.uk



## Monitoring the Grid is a Challenge

Number of participating sites is growing every day:

August 2003 => 12 sites ;

October 2004 => 83 sites ; 8000 CPUs; 96 PB Disk

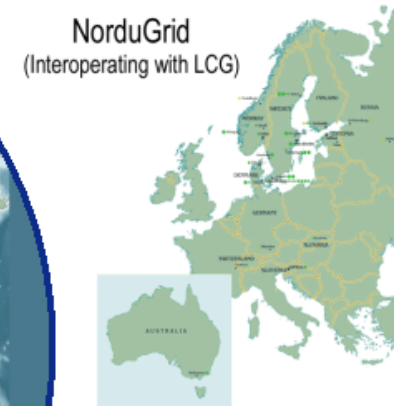
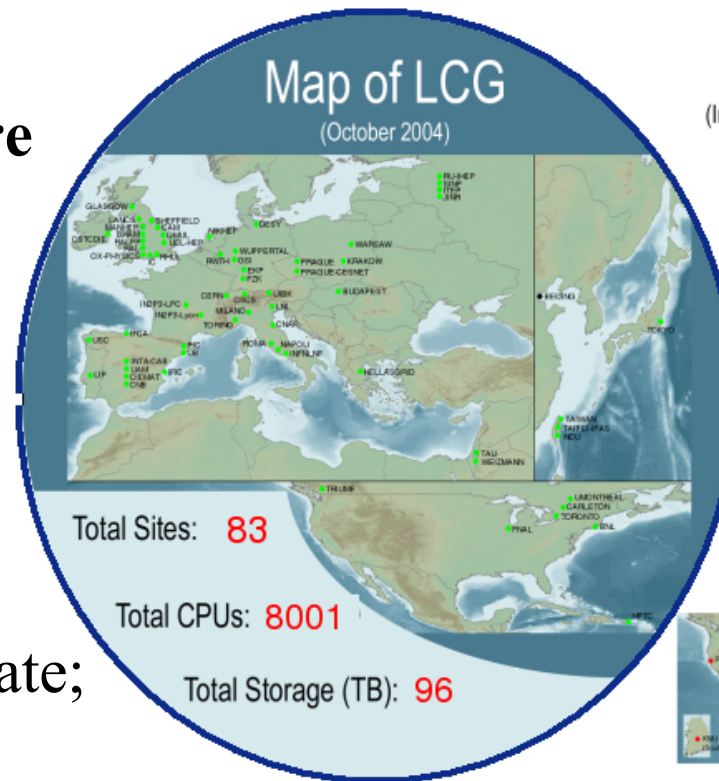
### Grid Operations Centre

Monitor the operational status of sites;

Fault detection

Problem Management

Identify problems; escalate; track;



## Introduction

- Look at the existing monitoring tools that are being used in LCG Grid Operations Centre
    - GPPMON
    - GRIDICE
    - GSTAT
    - CERTIFICATION TESTING
    - REAL TIME GRID MONITOR
    - Job Accounting
  - There has been a coordinated effort to develop, deploy and integrate a variety of monitoring tools from CERN, CCLRC (UK), GridPP, INFN-Grid (Italy) and Taiwan.
-

- We have only fragmentary information about the services that sites are running.
  - We don't know what RBs/SEs/Sites the VOs are using for data challenges.
  - We don't know what the core services are and who is running them.
  - We don't have a toolkit to test specific core services.
  - We have to concentrate on functional behaviour of services e.g. If an RB sends your job to a CE, then we must assume the RB is working fine. Is this the only test of a RB?
  - Not all the tests that we perform are effective at finding problems.
  - We must develop tests which simulate the life cycle of real applications in a Grid environment.
  - ...and lots more (see earlier talks)
-

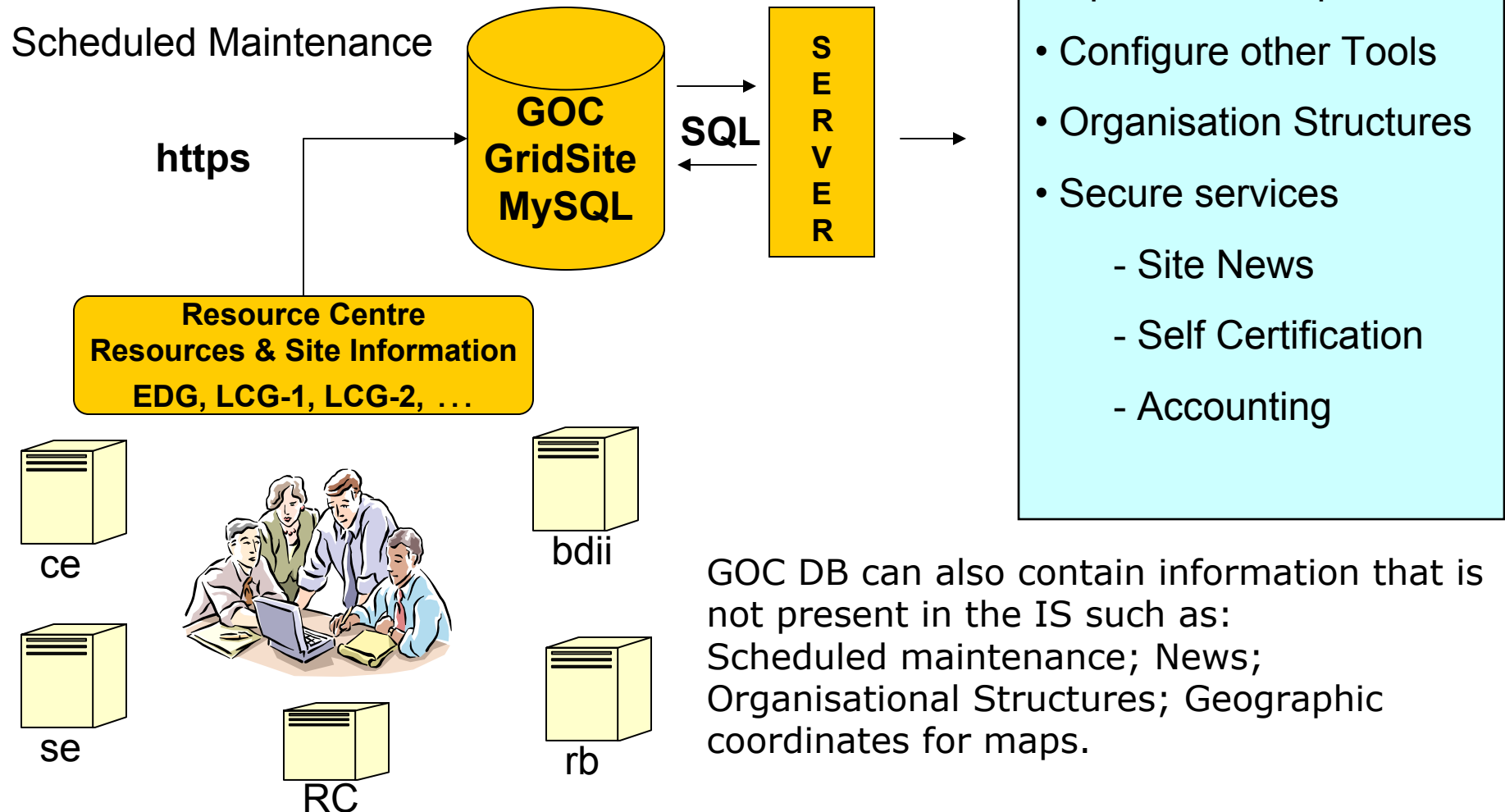
# CCLRC GOC Configuration Database

Secure Database Management via HTTPS / X.509

Store a Subset of the Grid Information system

People, Contact Information, Resources

Scheduled Maintenance





# Operations Map – Job Submission Tests

## GPPMON

Displays the results of tests against sites.

Test: Job Submission

Job is a simple test of the grid middleware components e.g. Gatekeeper service, RB service, and the Information System via JDL requirements.



This kind of test deals with the functional behaviour core grid services – do simple jobs run. They are lightweight tests which run hourly. However, they have certain limitations e.g. Dteam VO; WN reach (specialised monitoring queues).

# Operations Map – Certificate Lifetime

## GPPMON

Displays the results of tests against sites.

Test: Certificate Lifetime

Many grid services require a valid certificate for security.

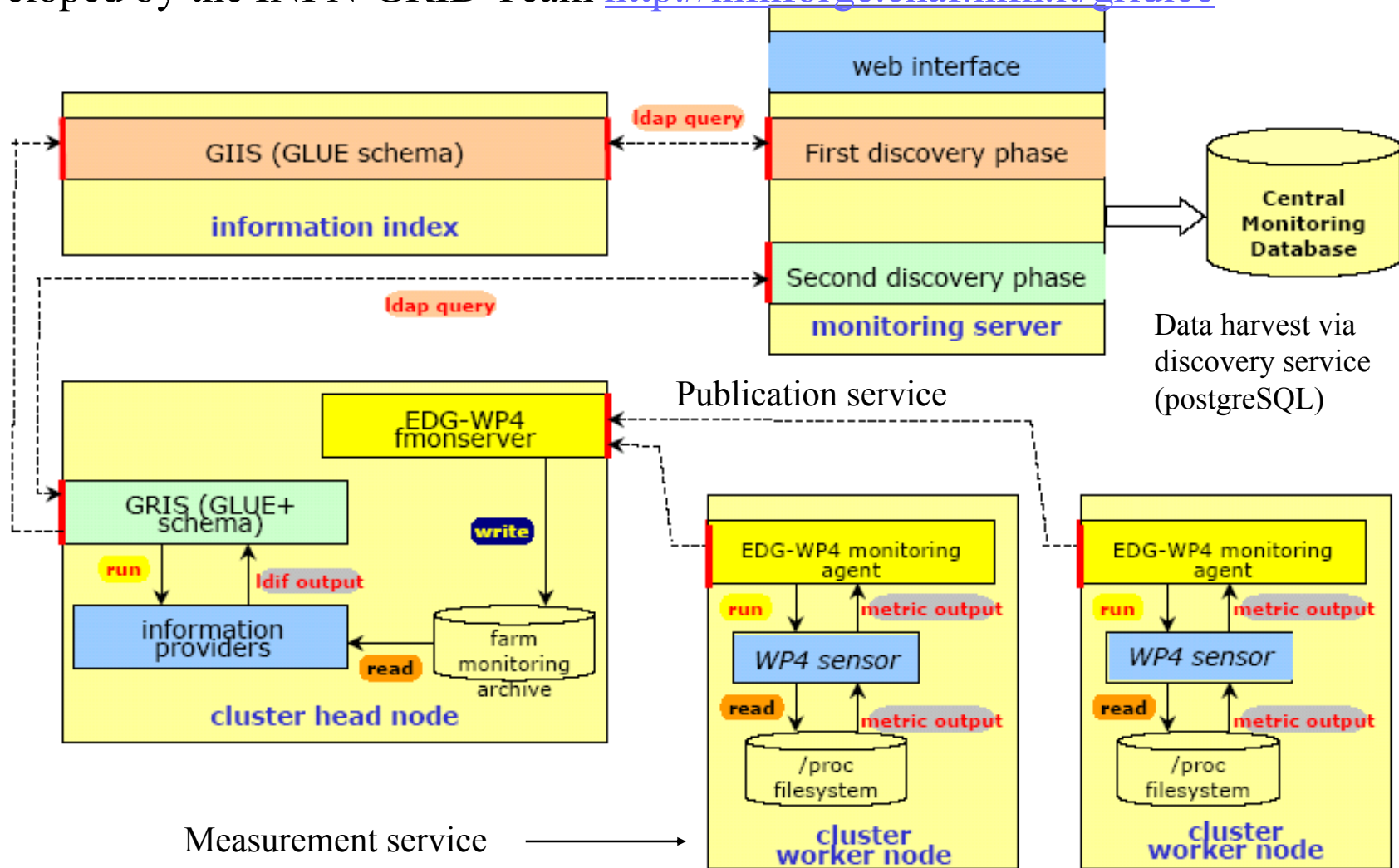


By probing the host certificates on CEs and SEs at sites with a simple SSL client service, we can identify certificates which are due to expire and send an early warning to them. A predictive tool!

# CCLRC GRIDICE – Architecture

A different kind of monitoring tool – processes / low level metrics / grid metrics

Developed by the INFN-GRID Team <http://infnforge.cnaf.infn.it/gridice>









# GRIDICE – Global View


Different Views of the data: Site / VO / Geographic



GOC



GridICE  
the eyes of the Grid



INFN  
GRID

Site view
VO view
Geo view
Gris view
Help
about

List of Sites

Select Site  and/or Role  Show

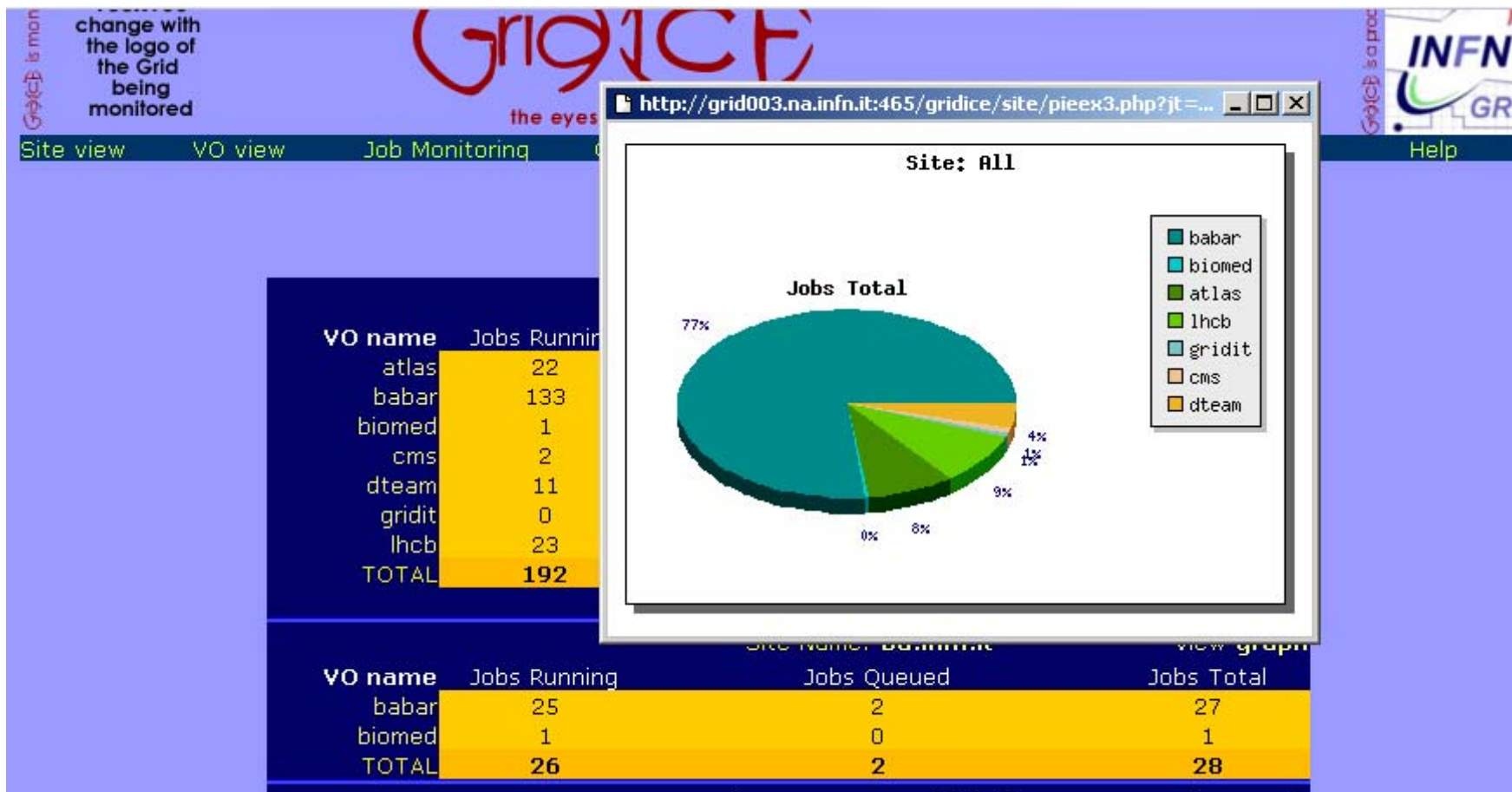
Site	Computing Resources					Storage Resources					
	Slot#	SlotFree	SlotLoad	RunJob	WaitJob	Power	CPU#	CPUload	Available	Total	%
cern.ch	408	180	55%	0	0	-	-	-	67.5 Gb	69.1 Gb	2%
cnaf.infn.it	-	-	-	-	-	-	-	-	-	-	-
cr.cnaf.infn.it	2154	1086	49%	253	0	762647	387	51%	868.0 Gb	999.7 Gb	13%
fnal.gov	12	12	0%	0	0	-	-	-	-	-	-
fzk.de	-	-	-	-	-	-	-	-	-	-	-
gridka.de	-	-	-	-	-	-	-	-	-	-	-
gridpp.rl.ac.uk	438	273	37%	55	0	-	-	-	59.8 Gb	69.0 Gb	13%
grid.sinica.edu.tw	294	294	0%	0	0	-	-	-	-	-	-
hep.ph.ic.ac.uk	126	126	0%	0	0	-	-	-	9.2 Gb	16.8 Gb	45%
ifae.es	480	480	0%	0	0	433978	160	0%	5.6 Tb	22.4 Tb	25%
nikhef.nl	500	230	54%	137	13	-	-	-	1.4 Tb	1.7 Tb	20%
triumf.ca	4490	30	99%	0	0	-	-	-	729.1 Gb	731.1 Gb	0%
<b>TOTAL</b>	<b>8902</b>	<b>2711</b>	<b>33%</b>	<b>445</b>	<b>13</b>	<b>1196625</b>	<b>547</b>	<b>26%</b>	<b>8.6 Tb</b>	<b>25.9 Tb</b>	<b>17%</b>

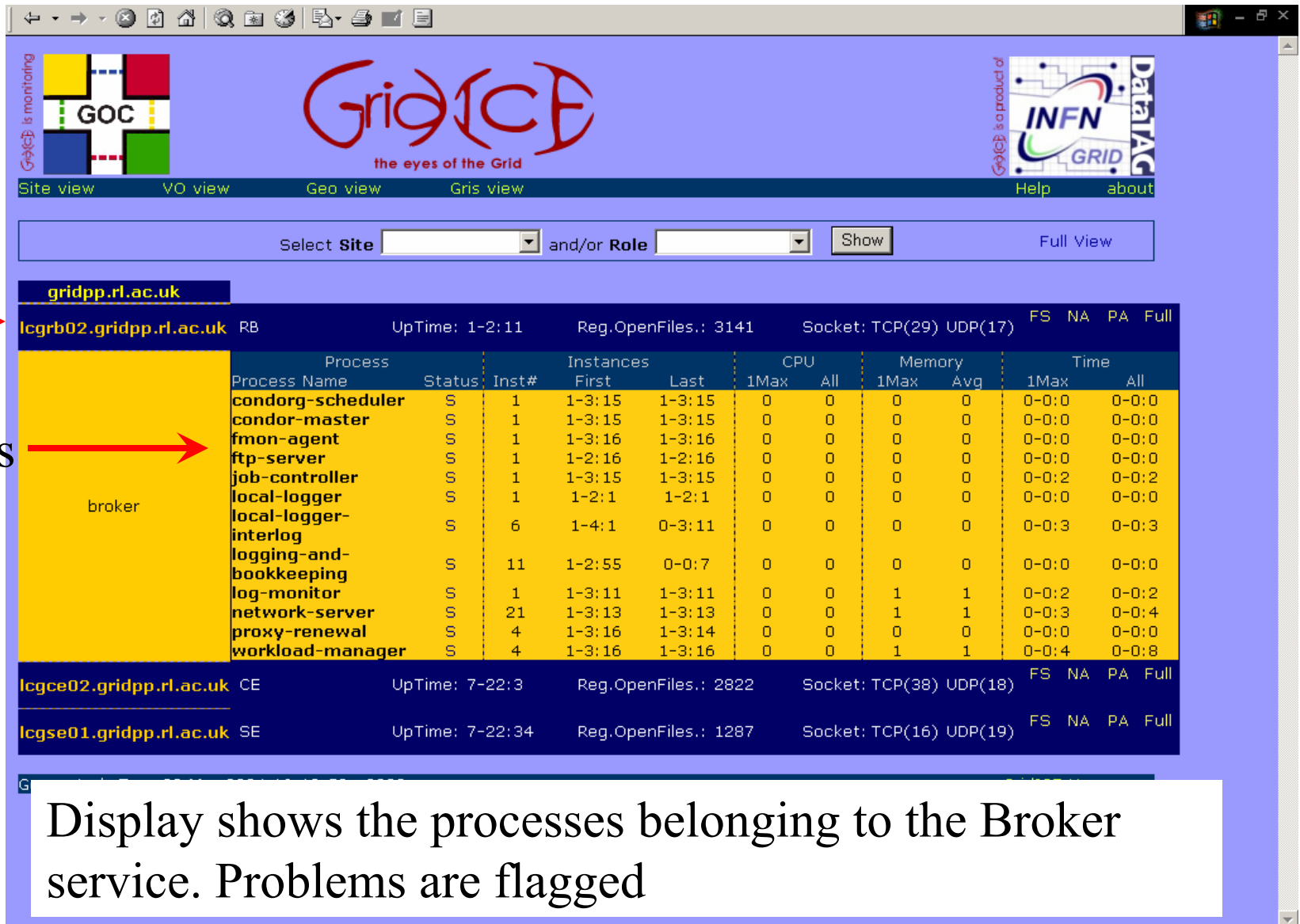
Generated: Mon, 22 Mar 2004 15:50:30 +0000
GridICE Homepage

SERVICES: PROGRAMS ARE TAGGED

# GridIce Job Monitoring

- Recently deployed version 1.6.3 on to LCG which features job monitoring: Queued, Running, Finished organised in different ways (site, Vo etc)
- XML views of data





The screenshot shows the GridICE Expert View interface. At the top, there are navigation tabs: Site view, VO view, Geo view, and Gris view. Below these are search filters for Site and Role, and a 'Show' button. The main content area displays a list of nodes under the heading 'gridpp.rl.ac.uk'. The selected node is 'lcgrb02.gridpp.rl.ac.uk', which is highlighted in yellow. A table of processes is shown for this node, with columns for Process Name, Status, Inst#, Instances (First, Last), CPU (1Max, All), Memory (1Max, Avg), and Time (1Max, All). The 'broker' process is highlighted in yellow. Below the table, there are summary statistics for other nodes: 'lcgce02.gridpp.rl.ac.uk' (CE) and 'lcgse01.gridpp.rl.ac.uk' (SE).

**Node** →

**Processes** →

Process Name	Status	Inst#	Instances First	Instances Last	CPU 1Max	CPU All	Memory 1Max	Memory Avg	Time 1Max	Time All
condorg-scheduler	S	1	1-3:15	1-3:15	0	0	0	0	0-0:0	0-0:0
condor-master	S	1	1-3:15	1-3:15	0	0	0	0	0-0:0	0-0:0
fmon-agent	S	1	1-3:16	1-3:16	0	0	0	0	0-0:0	0-0:0
ftp-server	S	1	1-2:16	1-2:16	0	0	0	0	0-0:0	0-0:0
job-controller	S	1	1-3:15	1-3:15	0	0	0	0	0-0:2	0-0:2
local-logger	S	1	1-2:1	1-2:1	0	0	0	0	0-0:0	0-0:0
local-logger-interlog	S	6	1-4:1	0-3:11	0	0	0	0	0-0:3	0-0:3
logging-and-bookkeeping	S	11	1-2:55	0-0:7	0	0	0	0	0-0:0	0-0:0
log-monitor	S	1	1-3:11	1-3:11	0	0	1	1	0-0:2	0-0:2
network-server	S	21	1-3:13	1-3:13	0	0	1	1	0-0:3	0-0:4
proxy-renewal	S	4	1-3:16	1-3:14	0	0	0	0	0-0:0	0-0:0
workload-manager	S	4	1-3:16	1-3:16	0	0	1	1	0-0:4	0-0:8

Display shows the processes belonging to the Broker service. Problems are flagged



# Ganglia Monitoring

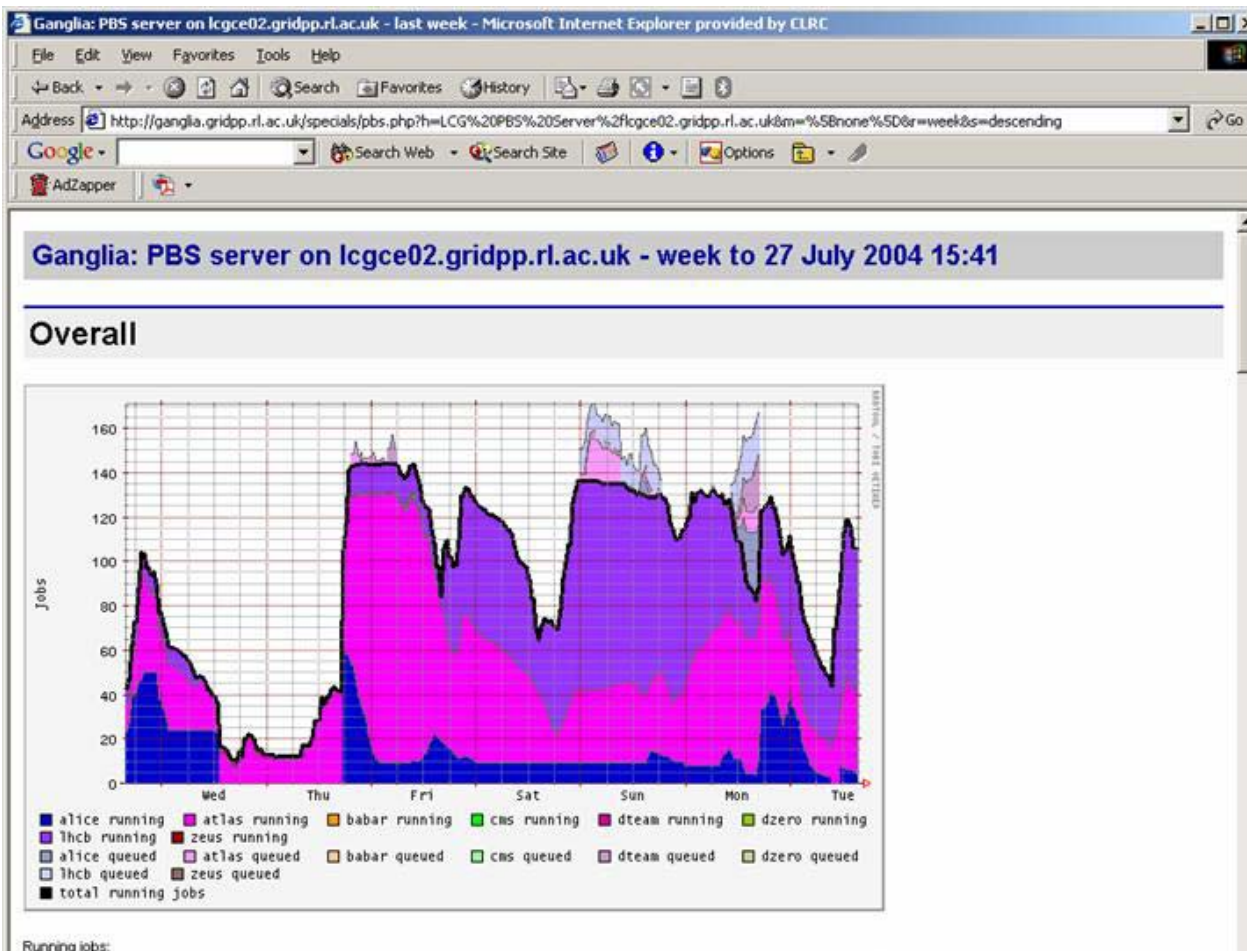
- <http://gridpp.ac.uk/ganglia>
- Can use Ganglia to monitor a cluster

**Scalable distributed monitoring system for clusters and grids using RRD for storage and visualisation.**

**RAL Tier-1 Centre**

**LCG PBS Server displays Job status for each VO**

**Get a lot for little effort**

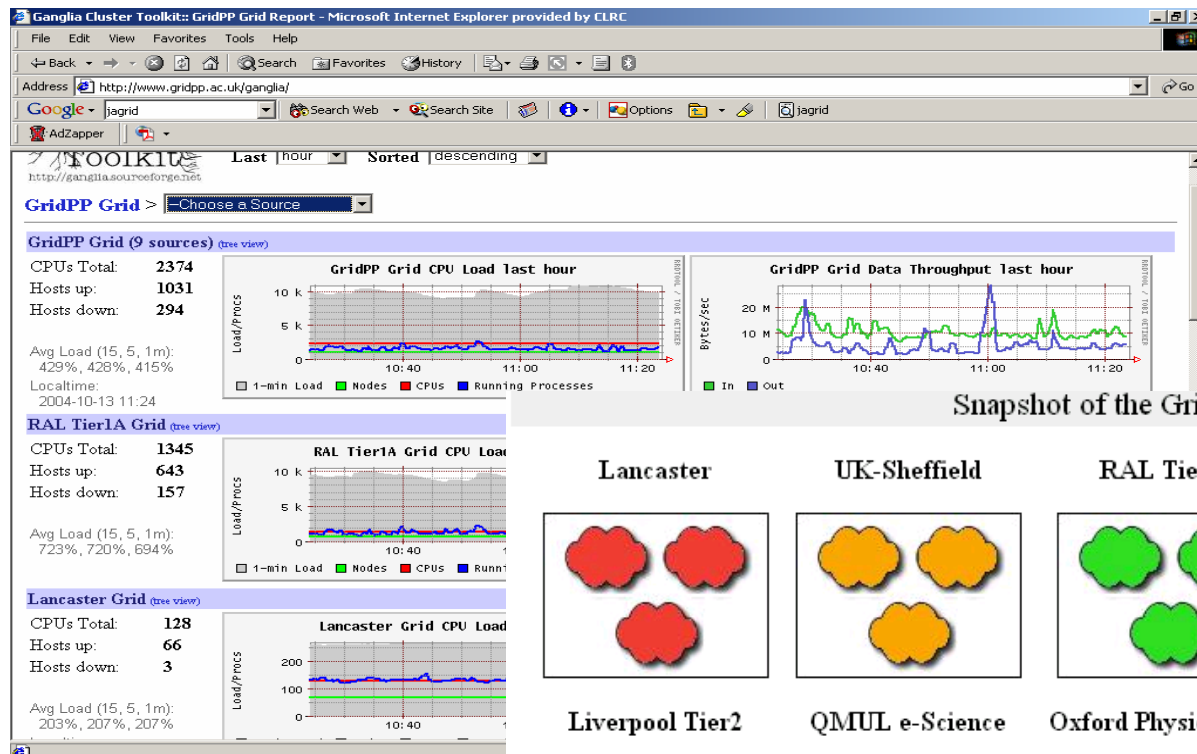




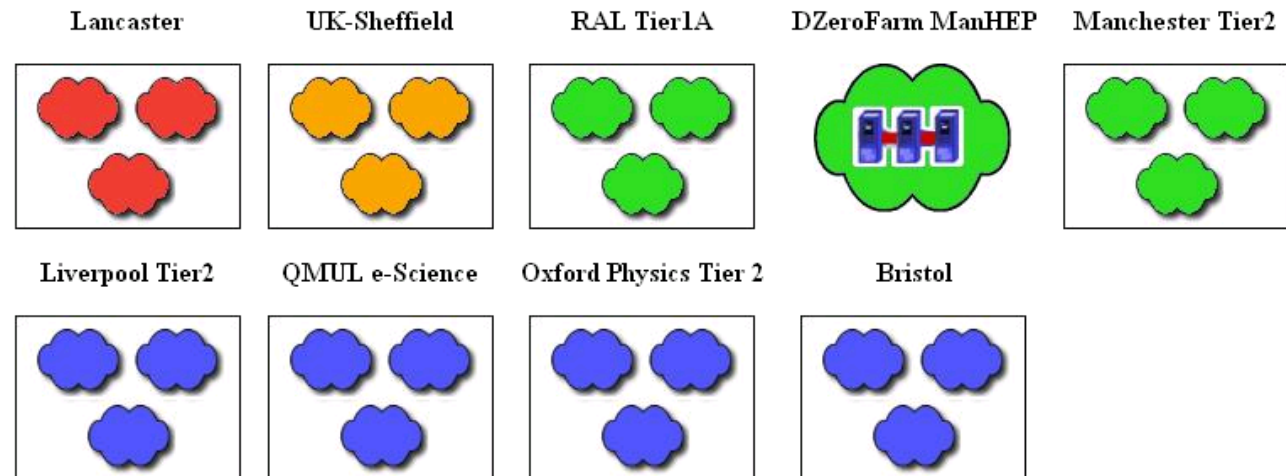
# Federating Cluster Information

- Can also use Ganglia to monitor clusters of clusters

Ganglia/R-GMA integration through Rangleia.



Snapshot of the GridPP Grid | Legend



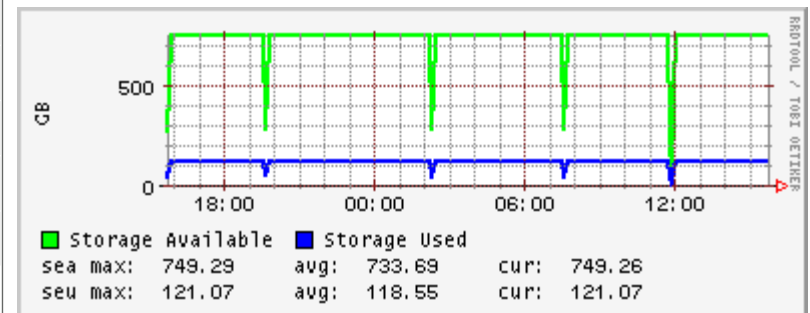
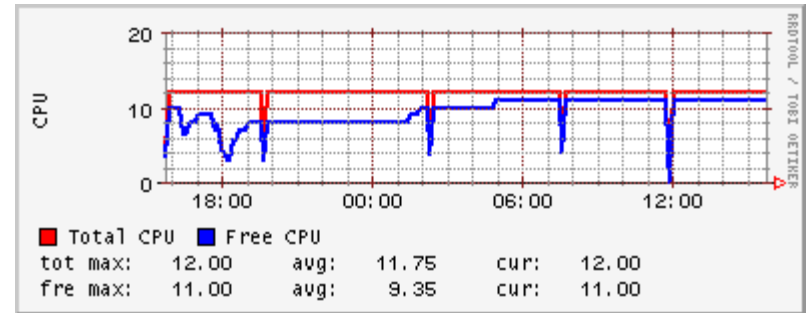
- Developed by MinTsai (GOC Taipei)
- Tool to display and check information published by the site GIIS (sanity checks, fault detection)
- <http://goc.grid.sinica.edu.tw/gstat/>

GIIS Monitor 15:43:32 10/21/04 GMT

[home](#) [table](#) [service](#) [regional](#) [help](#)

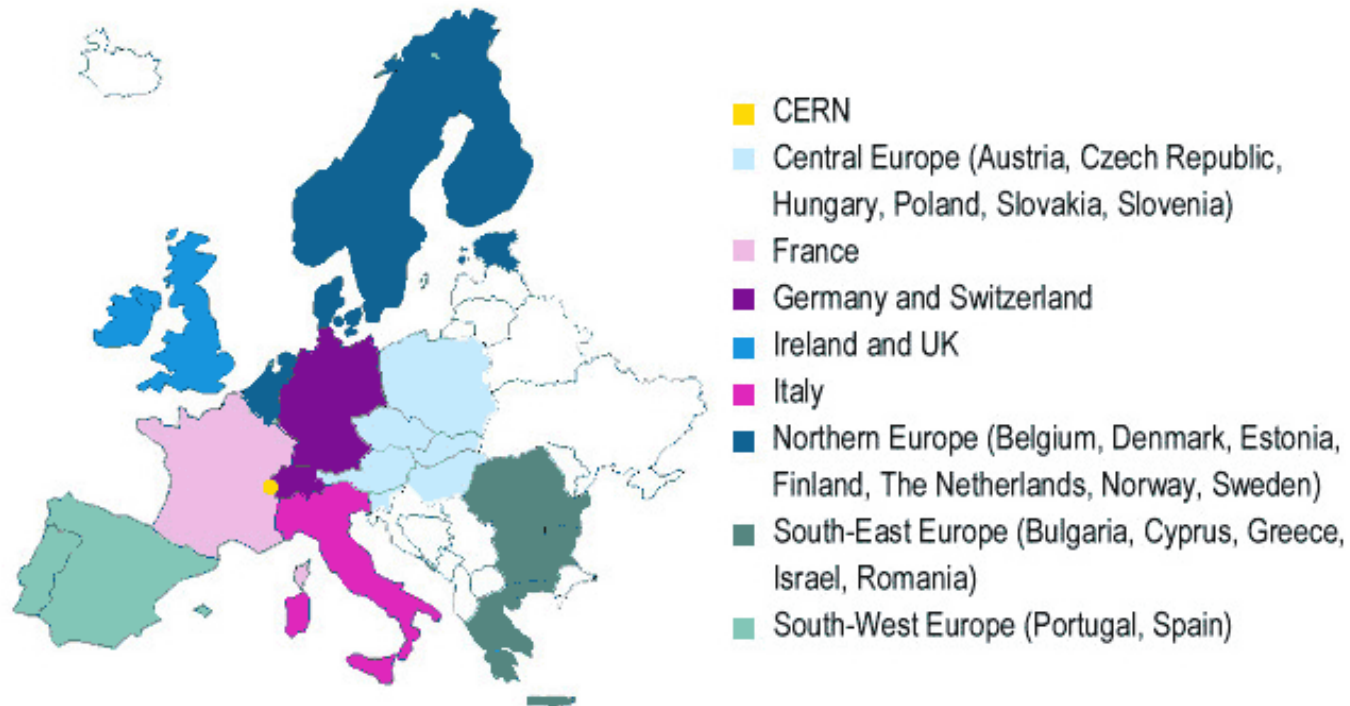
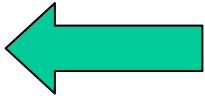
<a href="#">CAVENDISH-LCG2</a> <small>ok</small>	<a href="#">CERN-LCG2</a> <small>ok</small>	<a href="#">CIEMAT-LCG2</a> <small>ok</small>	<a href="#">CNAF-LCG2</a> <small>sw</small>	<a href="#">CYFRONET-LCG2</a> <small>ok</small>	<a href="#">FNAL-LCG2</a> <small>na</small>
<a href="#">IC-LCG2</a> <small>ok</small>	<a href="#">IFCA-LCG2</a> <small>ok</small>	<a href="#">IFIC-LCG2</a> <small>na</small>	<a href="#">INFN-LNL-LCG</a> <small>ok</small>	<a href="#">INFN-MILANO-LCG2</a> <small>ok</small>	<a href="#">INFN-TORINO-LCG2</a> <small>ok</small>
<a href="#">nikhef.nl</a> <small>ok</small>	<a href="#">PIC-LCG2</a> <small>ok</small>	<a href="#">RAL-LCG2</a> <small>sw</small>	<a href="#">Taiwan-LCG2</a> <small>sw</small>	<a href="#">Taiwan-NCU-LCG2</a> <small>ok</small>	<a href="#">Uni-Wuppertal</a> <small>ok</small>
<a href="#">USC-LCG2</a> <small>ok</small>	<a href="#">alberta-lcg2</a> <small>sw</small>	<a href="#">BEIJING-IHEP-LCG2</a> <small>ok</small>	<a href="#">BG01-IPP</a> <small>sw</small>	<a href="#">BHAM-LCG2</a> <small>ok</small>	<a href="#">BitLab-LCG2</a> <small>na</small>
<a href="#">BNL</a> <small>sw</small>	<a href="#">BUDAPEST</a> <small>ok</small>	<a href="#">CARLETONU-LCG2</a> <small>sw</small>	<a href="#">CCIN2P3-LCG2</a> <small>ok</small>	<a href="#">CGG-LCG2</a> <small>na</small>	<a href="#">CNB-LCG2</a> <small>ok</small>
<a href="#">CSCS-LCG2</a> <small>ok</small>	<a href="#">CY01-LCG2</a> <small>ok</small>	<a href="#">DESYPRO</a> <small>na</small>	<a href="#">ekplcg2</a> <small>ok</small>	<a href="#">FZK-LCG2</a> <small>na</small>	<a href="#">GSI-LCG2</a> <small>na</small>
<a href="#">HEPHY-UIBK</a> <small>ok</small>	<a href="#">HG-01-GRNET</a> <small>ok</small>	<a href="#">HPTC-LCG2</a> <small>ok</small>	<a href="#">IN2P3-LPC</a> <small>na</small>	<a href="#">INFN-FRASCATI</a> <small>na</small>	<a href="#">INFN-NAPOLI-ATLAS</a> <small>na</small>
<a href="#">INFN-ROMA1</a> <small>ok</small>	<a href="#">INTA-CAB</a> <small>ok</small>	<a href="#">IPSL-IPGP-LCG2</a> <small>ok</small>	<a href="#">ITEP</a> <small>ok</small>	<a href="#">JINR-LCG2</a> <small>ok</small>	<a href="#">LAL-LCG2</a> <small>na</small>
<a href="#">Lancs-LCG2</a> <small>ok</small>	<a href="#">LIP-LCG2</a> <small>na</small>	<a href="#">LivHEP-LCG2</a> <small>ok</small>	<a href="#">ManHEP-LCG2</a> <small>ok</small>	<a href="#">NCP-Lcg2</a> <small>na</small>	<a href="#">OXFORD-01-LCG2</a> <small>ok</small>
<a href="#">Prague-CESNET</a> <small>ok</small>	<a href="#">Prague-LCG2</a> <small>ok</small>	<a href="#">QMUL-eScience</a> <small>ok</small>	<a href="#">RALPP-LCG</a> <small>ok</small>	<a href="#">RHUL-LCG2</a> <small>ok</small>	<a href="#">ru-Moscow-KIAM-LCG2</a> <small>sw</small>
<a href="#">ru-Moscow-SINP-LCG2</a> <small>ok</small>	<a href="#">RU-Protvino-IHEP</a> <small>ok</small>	<a href="#">RWTH-LCG2</a> <small>sw</small>	<a href="#">SARA-LCG2</a> <small>ok</small>	<a href="#">ScotGRID-Edinburgh</a> <small>sw</small>	<a href="#">scotgrid-gla</a> <small>na</small>
<a href="#">SHEFFIELD-LCG2</a> <small>na</small>	<a href="#">Taiwan-IPAS-LCG2</a> <small>ok</small>	<a href="#">TAU-LCG2</a> <small>sw</small>	<a href="#">tiflog2</a> <small>na</small>	<a href="#">Tokyo-LCG2</a> <small>ok</small>	<a href="#">TORONTO-LCG2</a> <small>na</small>
<a href="#">TRIUMF-GC-CG2</a> <small>na</small>	<a href="#">TRIUMF-LCG2</a> <small>ok</small>	<a href="#">UAM-LCG2</a> <small>ok</small>	<a href="#">UB-LCG2</a> <small>na</small>	<a href="#">UCL-CCC</a> <small>ok</small>	<a href="#">UCL-HEP</a> <small>ok</small>
<a href="#">Umontreal-LCG2</a> <small>na</small>	<a href="#">WARSAW-LCG2</a> <small>ok</small>	<a href="#">WEIZMANN-LCG2</a> <small>sw</small>	<a href="#">(new)GR-01-AUTH</a> <small>ok</small>		

	sites	countries	totalCPU	freeCPU	runJob	waitJob	seAvail TB	seUsed TB	maxCPU	avgCPU
<b>Total</b>	82	26	7849	4073	2301	2305	96910.62	95024.78	16404	7975



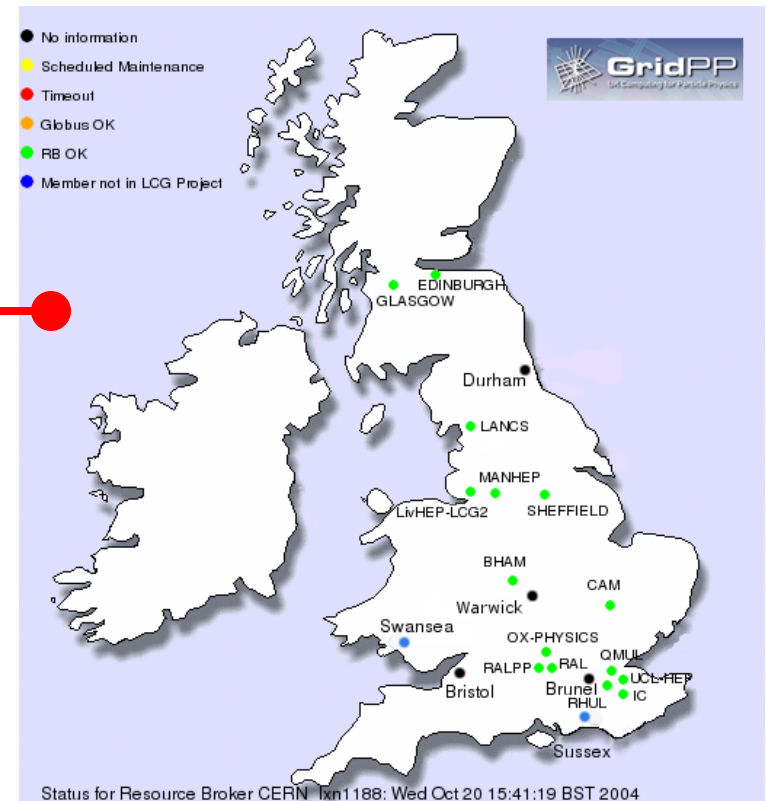
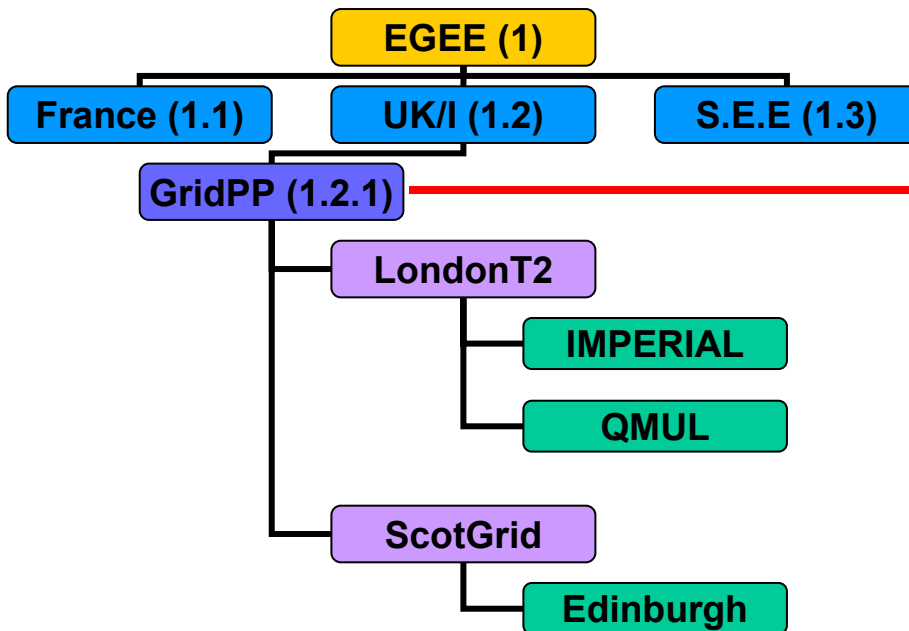
# Regional Monitoring

USA



- EGEE is made up of regions.
  - Each region contains many computing centres.
  - Regional Operational Centres is a focus for operations.
-

- ❑ [http://goc.grid-support.ac.uk/roc\\_map/map.php](http://goc.grid-support.ac.uk/roc_map/map.php)
- ❑ Provide ROCs with a package to monitor the resources in the region
  - Tailored Monitoring
  - GUIs to create organisations and populate them with sites
- ❑ Hierarchical view of Resources
  - Example UK Particle Physics GridPP
  - Materialised Path encoding





## Site Certification Service

- In terms of middleware, the installation and configuration of a site is quite a complicated procedure.
    - When there is a new release, sites don't upgrade at the same time
    - Some upgrades don't always go smoothly
    - Unexpected things happen (who turned off the power?)
    - Day-to-day problems; robustness of service under load?
  - It's necessary to actively hunt for problems
  - Site certification testing is by CERN deployment team on a daily basis. First step toward providing this service involves running a series of replica manager tests which register files onto the grid, move them around, delete them; and 3<sup>rd</sup> party copies from remote SE.
  - Unlike the simple job submission tests implemented in GPPMON, these tests are more heavy weight and attempt to simulate the life cycle of real applications.
-

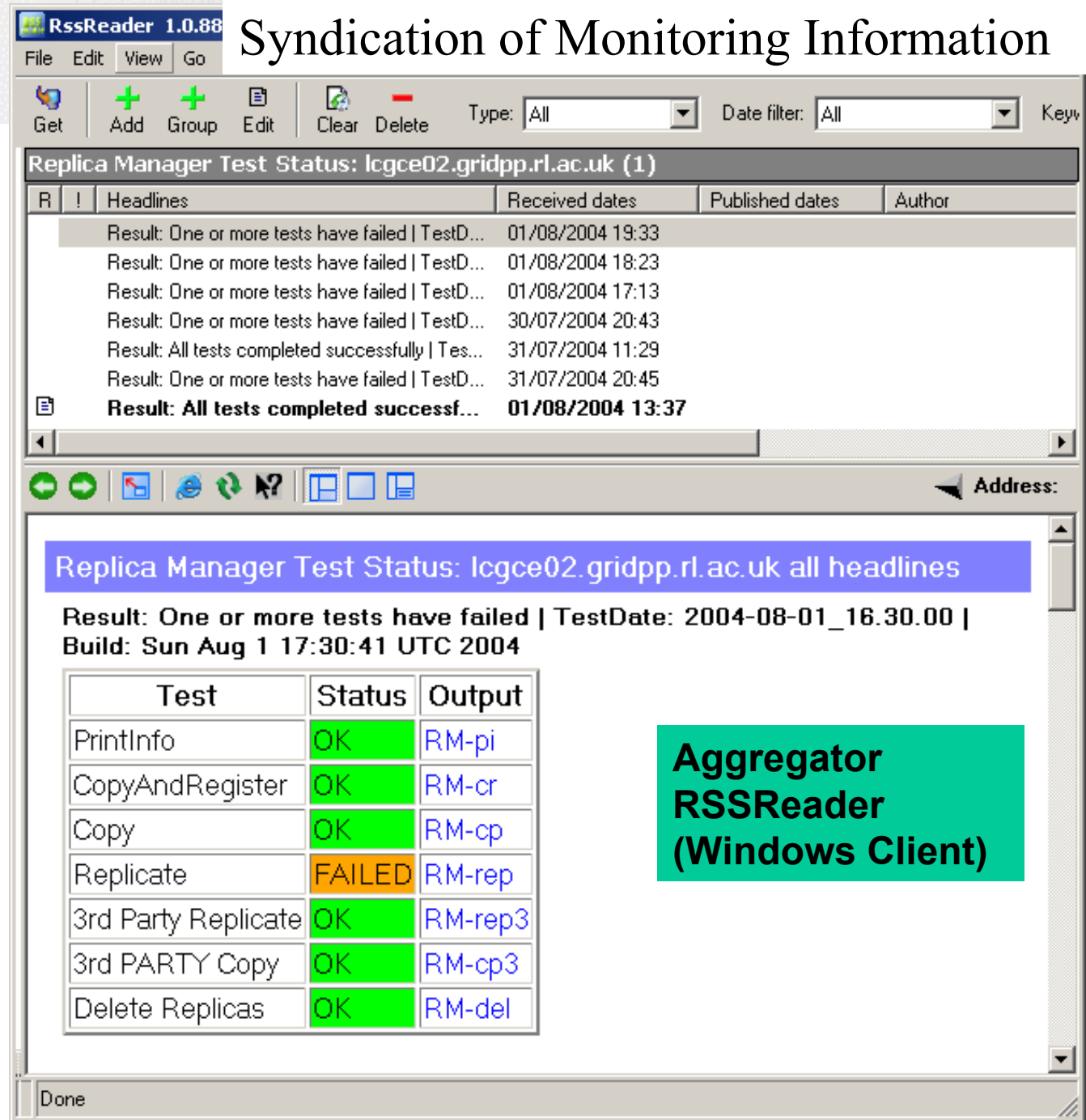


# RSS

GOC generates RSS feeds which clients can pull using an RSS aggregator.

How can we integrate feeds and ticketing systems?

## Syndication of Monitoring Information



The screenshot shows the RSSReader application window. The title bar reads "RssReader 1.0.88". The menu bar includes "File", "Edit", "View", and "Go". The toolbar contains icons for "Get", "Add", "Group", "Edit", "Clear", and "Delete", along with "Type:" and "Date filter:" dropdown menus. The main content area displays a list of headlines for "Replica Manager Test Status: lcge02.gridpp.rl.ac.uk (1)".

R	Headlines	Received dates	Published dates	Author
	Result: One or more tests have failed   TestD...	01/08/2004 19:33		
	Result: One or more tests have failed   TestD...	01/08/2004 18:23		
	Result: One or more tests have failed   TestD...	01/08/2004 17:13		
	Result: One or more tests have failed   TestD...	30/07/2004 20:43		
	Result: All tests completed successfully   Tes...	31/07/2004 11:29		
	Result: One or more tests have failed   TestD...	31/07/2004 20:45		
	<b>Result: All tests completed successf...</b>	<b>01/08/2004 13:37</b>		

The detailed view shows the following information:

**Replica Manager Test Status: lcge02.gridpp.rl.ac.uk all headlines**

**Result: One or more tests have failed | TestDate: 2004-08-01\_16.30.00 | Build: Sun Aug 1 17:30:41 UTC 2004**

Test	Status	Output
PrintInfo	OK	RM-pi
CopyAndRegister	OK	RM-cr
Copy	OK	RM-cp
Replicate	FAILED	RM-rep
3rd Party Replicate	OK	RM-rep3
3rd PARTY Copy	OK	RM-cp3
Delete Replicas	OK	RM-del

The status bar at the bottom of the window shows "Done".

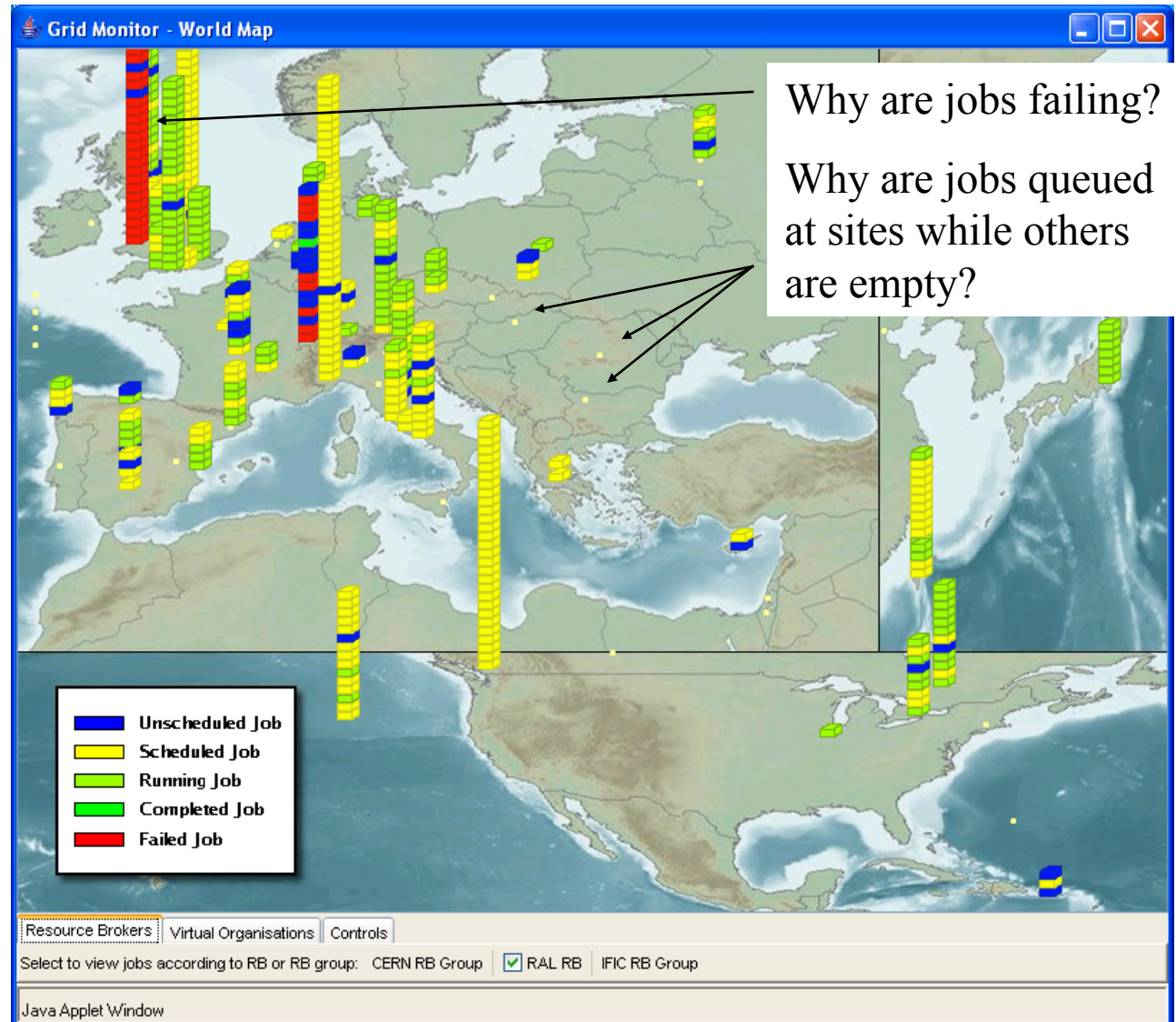
**Aggregator  
RSSReader  
(Windows Client)**

# Real Time Grid Monitor

<http://www.hep.ph.ic.ac.uk/e-science/projects/demo/index.html>

A Visualisation tool to track jobs currently running on the grid.

Applet queries the logging and bookkeeping service to get information about grid jobs.



# Problems with existing tools

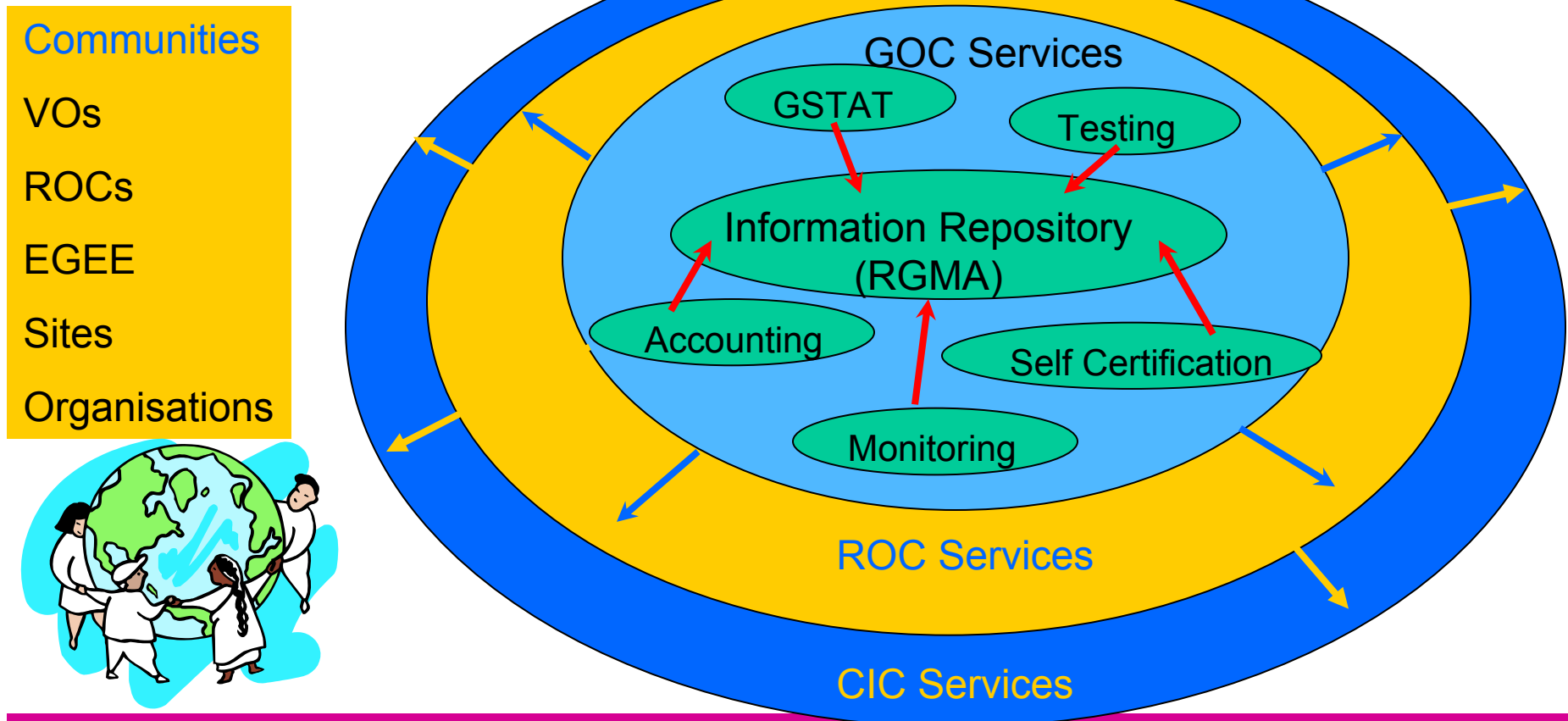
- Lots of monitoring tools have been described – they have a few things in common:
    - all the information which they generate is hidden away or difficult to access
    - limited interfaces: the data can only be accessed in specific ways
  - Therefore, it's difficult to build "on-demand" services to allow communities "Players" to interact with the data.
  - Examples include
    - a) Job Accounting service : to allow an Organisation to compare resources usage for each VO
    - b) Certification Testing service: Secure service to allow a site administrator to run the certification test suite against their site through a RB of their choice?
  - The idea is for the services to collect information and put it into a common repository such as an RGMA Archiver. In this way, the information can be shared and accessible to all.
  - Services (EGEE parlance: ROC and CIC services) munch the data and present it to the community.
  - Example: GIIS is that it's hard to drill down to the information you want e.g How much CPU in GridPP today? How much disk in the UKI ROC? The new paradigm solves this problem by allowing the data to be aggregated in different ways.
-

# Monitoring Paradigm

A Better way to unify monitoring information.

GOC Services collect information and publish into an archiver.

ROC/CIC Services provide a means for the community to interact with this information on-demand. GOC provides services tailored to the requirements of the community.

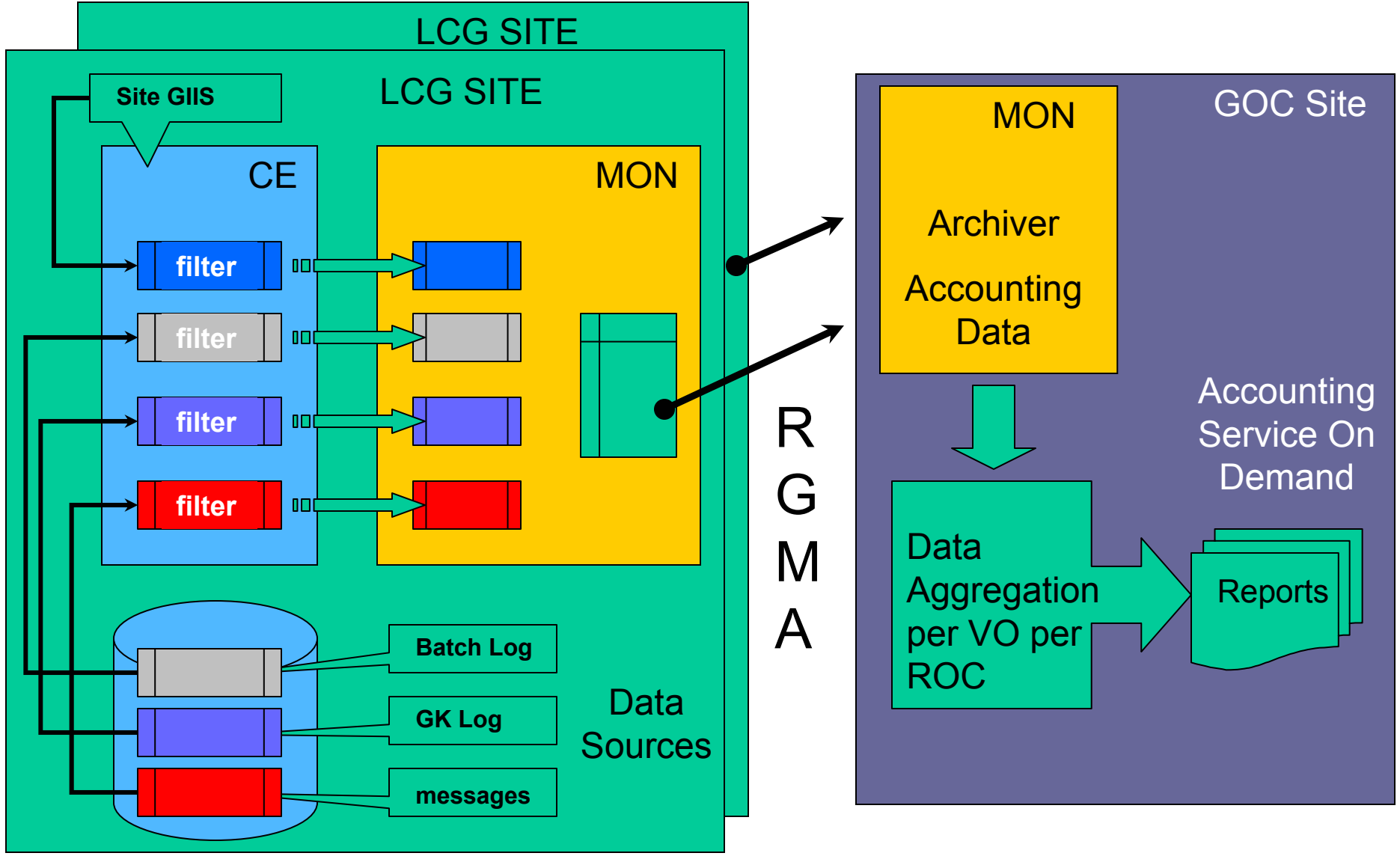




- An accounting package for LCG has been developed by the GOC at RAL
  - There are two main parts
    - the accounting data-gathering infrastructure based on R-GMA which brings the data to a central point
    - a web portal to allow on-demand reports for a variety of players.
-



# Accounting Flow Diagram





# GOC Accounting Services

GRID OPERATIONS CENTRE    News    Monitoring    Accounting    Operations    GOC Portal

**Accounting**

- Accounting Home
- CIC View
- ROC View
- Query Wizard
- Contact the GOC team

**CIC Accounting Services**

This view allows you to view summarised data across the entire EGEE organisation, sorted by VO.

NOTE: The range of VOs presented below are only those present in the test schema, and do not reflect the final list

Date range: Start year - 2004, Start month - 1, End year - 2004, End month - 4

VOs:  2688  alice  atlas  babar  cms  dteam  dzero  lhcb

Start point: EGEE

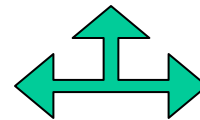
Refresh

Accounting plot:

On Demand Services to EGEE Community

Simple interface to customise views of data: VO, time frame and Region (default = EGEE)

BaseCpuSeconds Aggregated across EGEE



## Each Region, per VO, per Month

Group / Site	VO	Jan 2004	Feb 2004	Mar 2004	Apr 2004
Italy	2688	0	0	14619	0
	alice	0	5040	220333463	0
	atlas	0	3	877	0
	babar	0	0	0	0
	cms	0	0	5776	0
	dteam	0	1560	9208	0
	dzero	0	0	0	0
	lhcb	0	0	47	0

Group / Site	VO	Jan 2004	Feb 2004	Mar 2004	Apr 2004
UK / Ireland	2688	0	0	0	0
	alice	0	410303	94319599	589065
	atlas	4	0	1408	76006
	babar	25	0	7	3
	cms	0	47	1341	1137
	dteam	522	11719	24667	26962
	dzero	0	0	28	0
	lhcb	0	0	64	141

## Each Site, per VO, per Month

Group / Site	VO	Jan 2004	Feb 2004	Mar 2004	Apr 2004
CNAF	2688	0	0	14619	0
	alice	0	5040	220333463	0
	atlas	0	3	877	0
	babar	0	0	0	0
	cms	0	0	5776	0
	dteam	0	1560	9208	0
	dzero	0	0	0	0
	lhcb	0	0	47	0
	INFNLNF	2688	0	0	0
alice		0	0	0	0
atlas		0	0	0	0
babar		0	0	0	0
cms		0	0	0	0
dteam		0	0	0	0
dzero		0	0	0	0
lhcb		0	0	0	0
LEGNARO		2688	0	0	0
	alice	0	0	0	0
	atlas	0	0	0	0
	babar	0	0	0	0
	cms	0	0	0	0
	dteam	0	0	0	0
	dzero	0	0	0	0
	lhcb	0	0	0	0

Other Distributions  
Normalised CPU  
# Jobs

# Accounting Issues

1. A stable release of accounting package has been certified and tested at CERN; Should sites wait for the official release of press ahead independently?
  2. Package supports PBS only; initial implementation for LSF.  
80 sites advertising 313 Job managers:
    - 300 PBS (91% of sites)
    - 3 CONDOR (KFKI, FNAL, TRIUMF)
    - 7 LSF (GSI, LNL, CERN).
  3. Accounting requires the R-GMA infrastructure to be deployed at the site.
  4. The VO associated with a user's DN is not available in the batch or gatekeeper logs. It will be assumed that the group ID used to execute user jobs, which is available, is the same as the VO name.
  5. The global jobID assigned by the Resource Broker is not available in the batch or gatekeeper logs. This global jobID cannot therefore appear in the accounting reports. The RB Events Database contains this, but that is not accessible nor is it designed to be easily processed. [Andrea Guarise: JRA1 proposal]
-

# Accounting Issues

6. Most sites keep GK/Batch logs but throw away message log files after 9 weeks due to default log rotation.
  7. At present the logs provide no means of distinguishing sub-clusters of a CE which have nodes of differing processing power. Changes to the information logged by the batch system will be required before such heterogeneous sites can be accounted properly. At present it is believed all sites are homogeneous.
-

# Future Plans

- Extend the ideas developed in the accounting service to the other tools.
  - Example: Feeds
    - Regional Operations News feeds
    - (accounting, #cpu, disk, Piotrs Daily test results)
  - Want to move toward a Service Orientated Architecture model and provide the community with a direct interface into the monitoring.
-

# Summary

- Accounting Information gathering infrastructure has been developed
  - It has been through the C&T cycle and should be deployed in the next release.
  - A web portal for display of this information has been developed (work in progress)
  - This is an EGEE deliverable (DSA1.3)
  - The display infrastructure can be deployed for other monitoring information.
  - Development towards on-demand services to provide the community with up-to-date information, aggregated at different levels.
  - Development of Visualisation tools to enhance our understanding of the grid.
-

## Summary

- Since August 2003, the LCG GOC has been working to understand the problems of running a large scale distributed grid.
  - Setup a distributed GOC and deployed tools to help understand the issues.
  - Development towards on-demand services to provide the community with up-to-date information, aggregated at different levels.
  - Development of Visualisation tools to enhance our understanding of the grid.
-