

# LCG-1 Status

Markus Schulz  
LCG

EDG Project Conference  
29 September 2003



# Overview



- Our goals, scale, milestones (no visions etc.)
- Deployment status
- Software is in LCG-1 now
- Release and Deployment Procedures
- Services, Operation etc.
- First experience
- What do we plan to do in the near future (2003, mid 2004)
- Summary

Many slides stolen/inspired





# What is LCG?



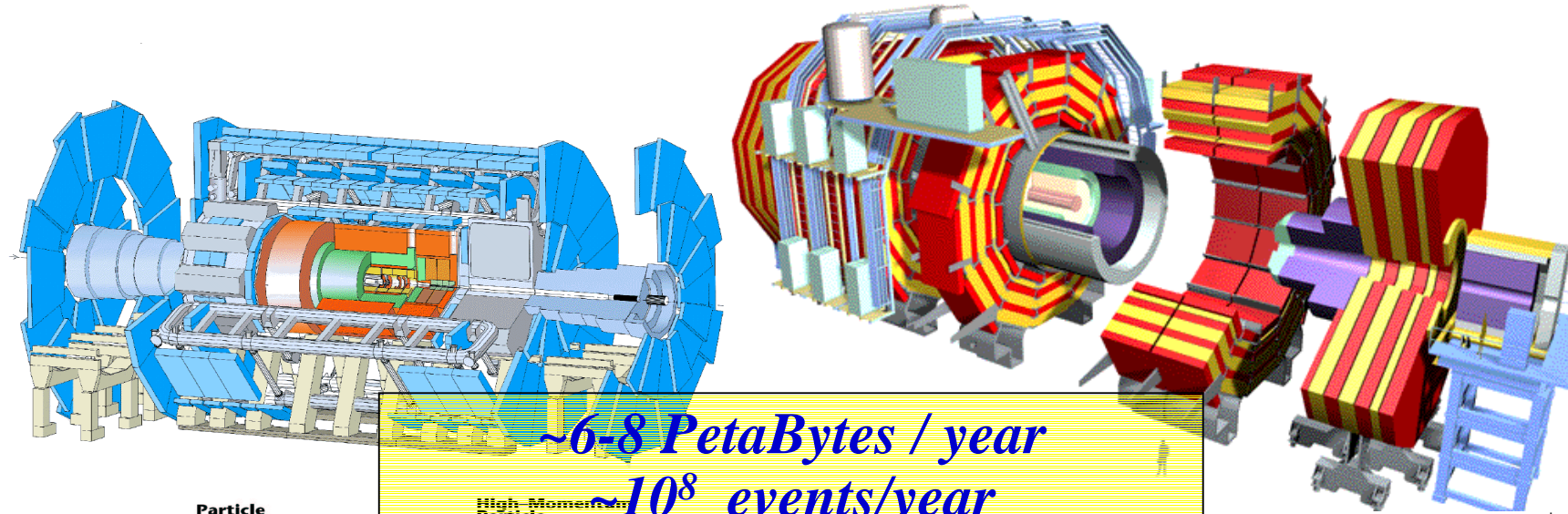
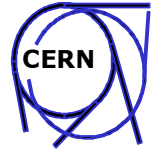
- **LHC Computing Grid** -> <http://lcg.web.cern.ch/lcg>
- The goal of the LCG project is to **prototype** and **deploy** the computing environment for the LHC experiments
- Two phases:
  - **Phase 1: 2002 – 2005**
  - Build a service prototype, based on existing grid middleware
  - Gain experience in running a **production grid service**
  - Produce the TDR for the final system
  - **Phase 2: 2006 – 2008**
  - Build and commission the initial LHC computing environment

**LCG is NOT a development project**

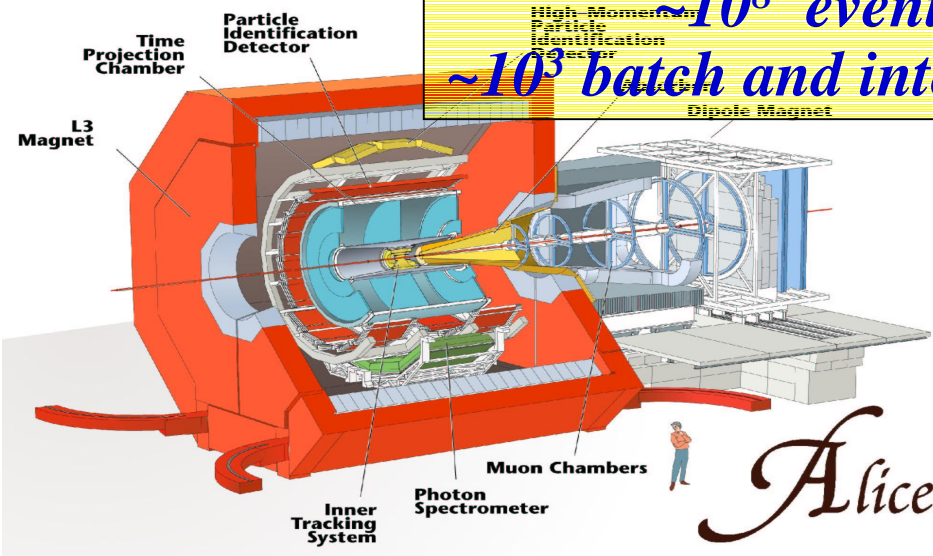




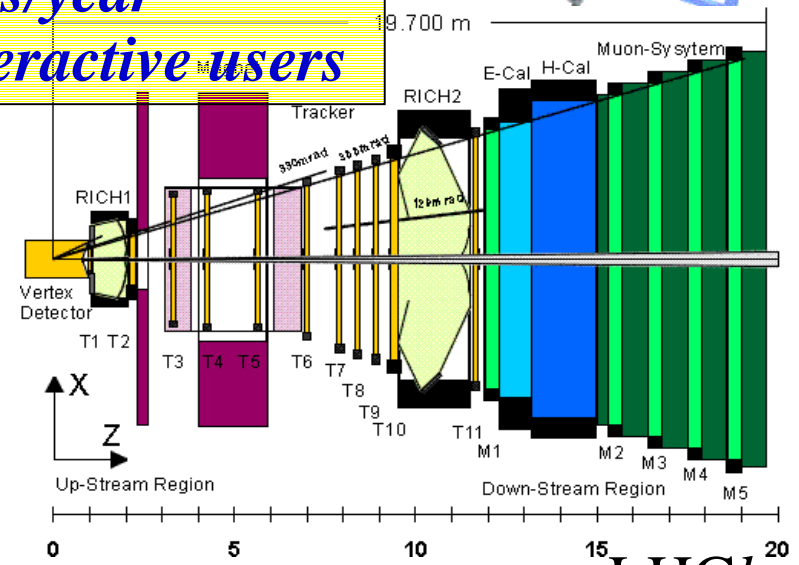
# Our Customers



~6-8 PetaBytes / year  
 ~10<sup>8</sup> events/year  
 ~10<sup>3</sup> batch and interactive users



*Alice*

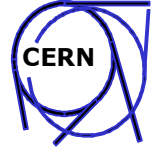


Federico.carminati , EU review presentation

LHCb



# 2003 Milestones

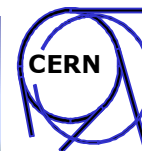


## Project Level 1 Deployment milestones that had been set for 2003:

- July: Introduce the initial publicly available LCG-1 global grid service
  - With 10 Tier 1 centres in 3 continents
- November: Expanded LCG-1 service with resources and functionality sufficient for the 2004 Computing Data Challenges
  - Additional Tier 1 centres, several Tier 2 centres – more countries
  - Expanded resources at Tier 1s
    - (e.g. at CERN make the LXBatch service grid-accessible)
  - Agreed performance and reliability targets



# LCG Resources (promised) – 1Q04



|                                      | <i>CPU<br/>(kSI2K)</i> | <i>Disk<br/>TB</i> | <i>Support<br/>FTE</i> | <i>Tape<br/>TB</i> |
|--------------------------------------|------------------------|--------------------|------------------------|--------------------|
| <i>CERN</i>                          | 700                    | 160                | 10.0                   | 1000               |
| <i>Czech Republic</i>                | 60                     | 5                  | 2.5                    | 5                  |
| <i>France</i>                        | 420                    | 81                 | 10.2                   | 540                |
| <i>Germany</i>                       | 207                    | 40                 | 9.0                    | 62                 |
| <i>Holland</i>                       | 124                    | 3                  | 4.0                    | 12                 |
| <i>Italy</i>                         | 507                    | 60                 | 16.0                   | 100                |
| <i>Japan</i>                         | 220                    | 45                 | 5.0                    | 100                |
| <i>Poland</i>                        | 86                     | 9                  | 5.0                    | 28                 |
| <i>Russia</i>                        | 120                    | 30                 | 10.0                   | 40                 |
| <i>Taiwan</i>                        | 220                    | 30                 | 4.0                    | 120                |
| <i>Spain</i>                         | 150                    | 30                 | 4.0                    | 100                |
| <i>Sweden</i>                        | 179                    | 40                 | 2.0                    | 40                 |
| <i>Switzerland</i>                   | 26                     | 5                  | 2.0                    | 40                 |
| <i>UK</i>                            | 1656                   | 226                | 17.3                   | 295                |
| <i>USA</i>                           | 801                    | 176                | 15.5                   | 1741               |
| <b>Total (1kSI2K is a 2.8GHz P4)</b> | <b>5600</b>            | <b>1169</b>        | <b>120.0</b>           | <b>4223</b>        |



# LCG-1 Deployment Status



- Up to date status can be seen here:
  - <http://www.grid-support.ac.uk/GOC/Monitoring/Dashboard/dashboard.html>
    - Has links to maps with sites that are in operation
    - Links to GridICE based monitoring tool (history of VO's jobs, etc)
      - Using information provided by the information system
    - Tables with deployment status
- Sites that are currently in LCG-1 ([here](#)) expect 18-20 by end of 2003
  - PIC-Barcelona (RB)
  - Budapest (RB)
  - CERN (RB)
  - CNAF (RB)
  - FermiLab. (FNAL)
  - FZK
  - Krakow
  - Moscow (RB)
  - RAL (RB)
  - Taipei (RB)
  - Tokyo

## Total number of CPUs ~120 WNs

### Sites to enter soon

BNL, Prague,(Lyon)

Several tier2 centres  
in Italy and Spain

### Sites preparing to join

Pakistan, Sofia,  
Switzerland

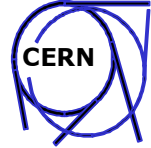
### Users (now):

Loose Cannons  
Deployment Team  
Experiments starting  
(Alice, ATLAS,..)

Some comments  
later....

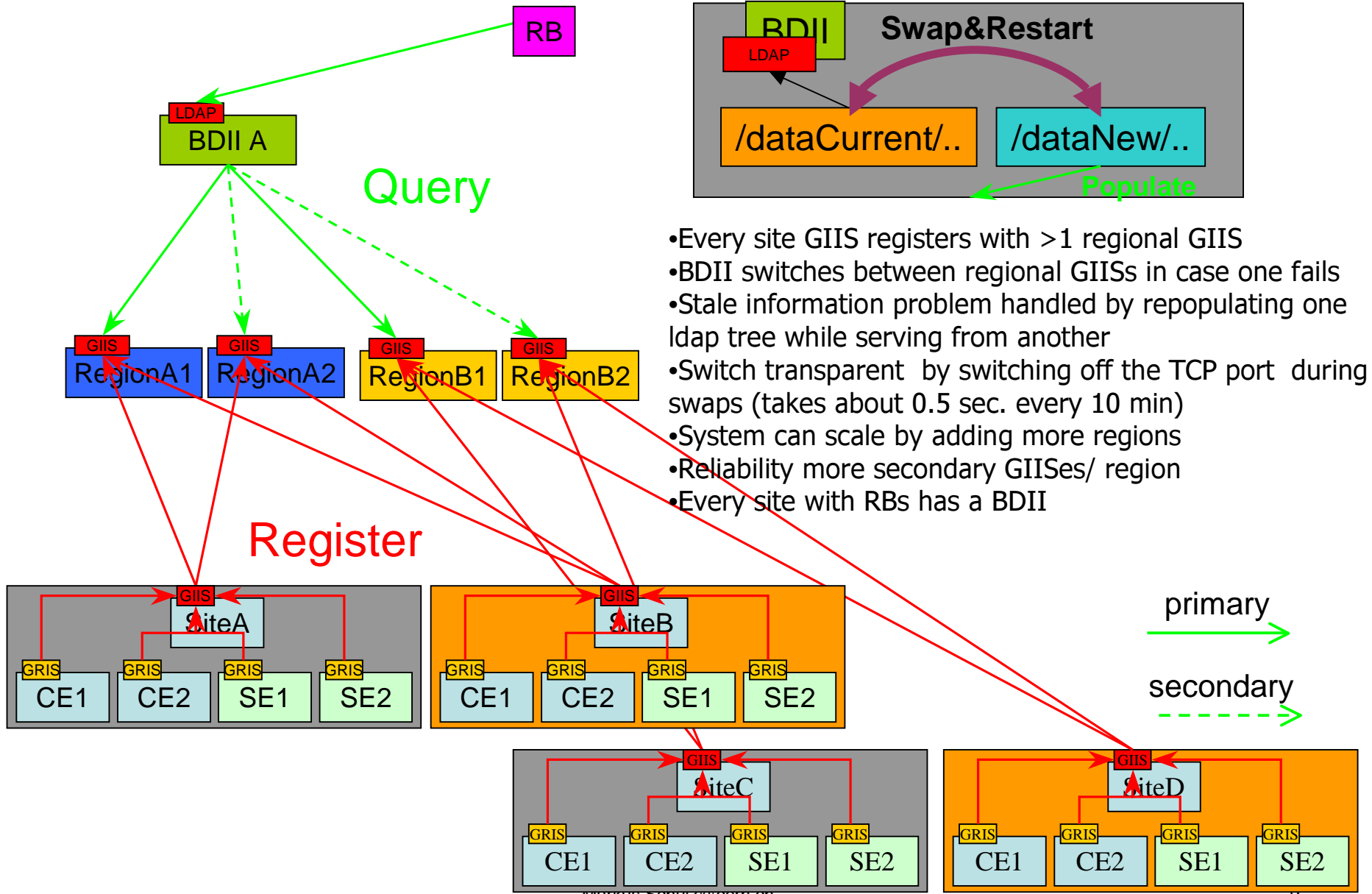


# LCG-1 Software



- LCG-1 (LCG1-1\_0\_2) is:
  - VDT (Globus 2.2.4)
  - EDG WP1 (Resource Broker)
  - EDG WP2 (Replica Management tools)
    - One central RMC and LRC for each VO, located at CERN, ORACLE backend
  - Several bits from other WPs (Config objects, InfoProviders, Packaging...)
  - GLUE 1.1 (Information schema) + few essential LCG extensions
  - MDS based Information System with LCG enhancements
  - SE-Classic (disk based only, gridFTP) **NO MSS**
  - EDG components approx. edg-2.0 version
  - LCG modifications:
    - Job managers to avoid shared filesystem problems (GASS Cache, etc.)
    - MDS – BDII LDAP (see more later)
    - Globus gatekeeper enhancements ((adding some accounting and auditing features, log rotation, that LCG requires)
    - Many, many bug fixes to EDG and Globus/VDT

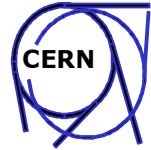




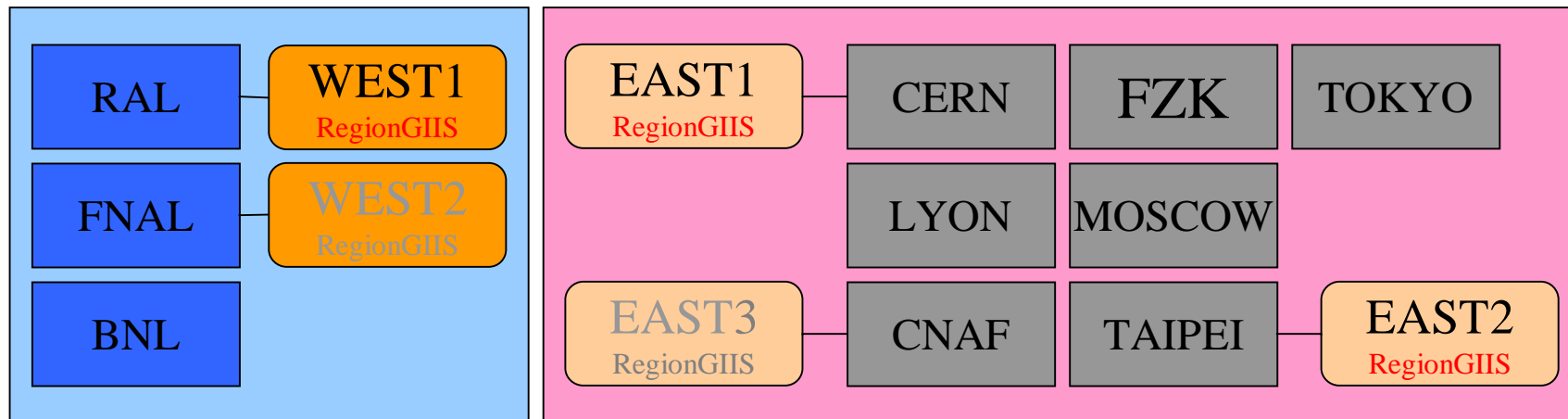
- Every site GIIS registers with > 1 regional GIIS
- BDII switches between regional GIISs in case one fails
- Stale information problem handled by repopulating one ldap tree while serving from another
- Switch transparent by switching off the TCP port during swaps (takes about 0.5 sec. every 10 min)
- System can scale by adding more regions
- Reliability more secondary GIISes/ region
- Every site with RBs has a BDII



# LCG-1 IS



- Current separation of the World:
  - Limit the number of sites/region
  - Started with 2, split along 0 degree into east and west
  - Currently 2 regional GIISes in East, only one in West (deployment)





# Release Procedure



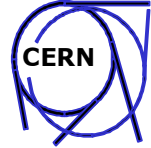
- Similar to what has been discussed for EDG in the past
  - Software first assembled on the Certification & Test Testbeds
    - 4 sites at CERN, + some external sites U. Wisconsin, FNAL, Moscow, Italy (soon)
    - Installation tests and functional test (resolving problems found in the LCG1 service)
    - Certification test suite almost finished
  - Software handed to the Deployment Team
    - Adjustments in the configuration
    - Release notes for the external sites
    - Decision on time to release
- How do we deploy?
  - Service Nodes (RB, CE, SE ...)
    - LCFGng, sample configurations in CVS
    - We provide for new sites config files based on a questionnaire
  - Worker nodes – aim is to allow sites to use existing tools as required
    - LCFGng – provides automated installation **YES**
    - Instructions allowing system managers to use their existing tools **SOON**
  - User interface
    - LCFGng **YES**
    - Installed on a cluster (e.g. Lxplus at CERN) LCFGng-lite **YES**
    - Instructions allowing system managers to use their existing tools **SOON**

Work intensive, limited to <10 sites





# Services



- **Operations Service:**
  - RAL is leading sub-project on developing operations services
  - Initial prototype <http://www.grid-support.ac.uk/GOC/>
    - Basic monitoring tools
    - Mail lists and rapid communications/coordination for problem resolution
    - Working on defining policies for operation, responsibilities (draft document)
  - **Monitoring:**
    - GridICE (development of DataTag Nagios-based tools)  
<http://tbed0116.cern.ch/gridice/site/site.php>
    - GridPP job submission monitoring <http://esc.dl.ac.uk/gppmonWorld/>
- **User support**
  - FZK leading sub-project to develop user support services
  - Draft on user support policy
  - Web portal for problem reporting <http://gus.fzk.de/>

# LCG Grid Operations Centre







Site View VO view about

| Site                | Info from Grid Discovery System |           |                   |              | Phys CPUs     |      | JOBS  |     |       |
|---------------------|---------------------------------|-----------|-------------------|--------------|---------------|------|-------|-----|-------|
|                     | Computing power                 | %JOB Load | %Storagearea Load | Tot. SA (GB) | SA Avail (GB) | Tot. | %Load | Run | Queue |
| cern.ch             | n/a                             | n/a       | 100%              | 87           | 78            | n/a  | n/a   | 0   | 0     |
| cr.cnaf.infn.it     | n/a                             | n/a       | 100%              | 63           | 56            | n/a  | n/a   | 0   | 0     |
| fbk.de              | n/a                             | n/a       | 100%              | 177          | 168           | n/a  | n/a   | 0   | 0     |
| grid.sinica.edu.tw  | n/a                             | n/a       | 26%               | 208          | 204           | n/a  | n/a   | 0   | 0     |
| icpp.su-tokyo.ac.jp | n/a                             | n/a       | 100%              | 63           | 60            | n/a  | n/a   | 0   | 0     |

- 📍 [Mapcenter](#) Monitoring of the LCG-1
- 📍 [Status Map](#) showing the status of nodes with respect to submitted test globus jobs
- 📍 [GRIDICE](#) Monitoring fo LCG-1 at CERN
- 📍 [Status table](#) of LCG-1 sites

| Site Name           | Country | Site Type | IPs | Resolution Class | Operation | Connected | Problems |
|---------------------|---------|-----------|-----|------------------|-----------|-----------|----------|
| cern.ch             | CH      | LCG-1     | 10  | LCG-1            | OK        | 100%      | 0        |
| cr.cnaf.infn.it     | IT      | LCG-1     | 10  | LCG-1            | OK        | 100%      | 0        |
| fbk.de              | DE      | LCG-1     | 10  | LCG-1            | OK        | 100%      | 0        |
| grid.sinica.edu.tw  | TW      | LCG-1     | 10  | LCG-1            | OK        | 100%      | 0        |
| icpp.su-tokyo.ac.jp | JP      | LCG-1     | 10  | LCG-1            | OK        | 100%      | 0        |



# Sites in LCG-1





# Global Grid User Support



- Current Status
- Documentation
- Training
- News
- Download
- Staff-internal

- Service Request
- Service Request
- the Knowledge Base
- the FAQs
- t

## News

**2003-09-08 11:18:12 UTC+0200**

In the Grid computing centre (GridKa) located in the research center Karlsruhe, the world ...

[...more](#)

**2003-09-08 11:18:43 UTC+0200**

Starting May 22nd, 2003 the Central GRID Support will launch a new version of the ...

[...more](#)

**2003-09-08 10:42:37 UTC+0200**

Due to urgent Hardware maintenance, the GRID center in KA will be down from 11:00 ...

[...more](#)

**2003-09-08 10:38:37 UTC+0200**

On June 12th, 2003 there will be an additional GRID user course for all ...

[...more](#)

## Status Information

**The current status of PBS GridKa at 2003-09-22 11:22:10 UTC+0200:**

The overall status is:



PBS Status is:

Active

Held Jobs:

0

Max configured Jobs are:

440

Total Jobs in Queues are:

1299

Running Jobs are:

439

Queued Jobs are:

859

Exiting Jobs are:

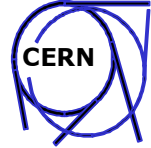
1

Wait Jobs in Queue:

0



# User Support

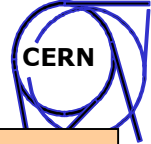


- Experiments provide 1st level triage
  - Experiments contacts send problems through the FZK portal
  - First test very soon, since the experiments start using LCG1
- Experiment integration support by CERN based group
- <http://grid-deployment.web.cern.ch/grid-deployment/cgi-bin/index.cgi?var=eis/homepage>
- Documentation
  - Installation guides (do first this, then that, if this happens do that..)
  - First version of user guide (very useful document)
  - Missing:
    - Collection of sample jobs
    - Tutorial
    - Operations manual





# Getting the Experiments on



No better place found for this slide

- Experiments start to use the service **now** and are welcome!!
  - Agreement between LCG and the experiments
    - System has limitations, testing what is there
    - Focus on:
      - Testing with loads similar to production programs (long jobs, etc)
      - Testing the experiments software on LCG
    - We don't want:
      - Destructive testing to explore the limits of the system with artificial loads
        - » This can be done in scheduled sessions on **C&T** testbed
  - Adding experiments and sites at a brisk pace in parallel is problematic
    - Getting the experiments on one after the other
      - **A** can learn from what **B** went through (**B** can claim fame for being 1st)
    - Limited number of users that we can interact with and keep informed
      - JJ's famous "Deadly Embrace until things are working"

Testing 10k "hello world" jobs on a system with 120CPUs doesn't help much to understand what has to be done to get 240 production jobs running for 12h.

**LCG needs the experiments  
NOW**



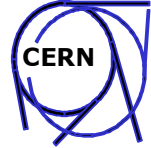
# Security



- LCG Security Group (Dave Kelsey (RAL))
  - LCG1 usage rules
  - Registration procedures and VO management
    - Agreement to collect only minimal amount of personal data
    - Currently registration is only valid for 6 month (procedures will change)
  - Initial audit requirements are defined
  - Initial incident response procedures
    - Site security contacts etc. are defined
  - Set of trusted CAs (including Fermilabs online KCA)
  - Draft of security policy (to be finished end of year)
  - Web site <http://proj-lcg-security.web.cern.ch/proj-lcg-security/>



# History

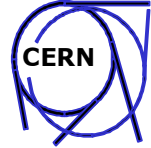


- First set of reasonable middleware on C&T Testbed end of July (PLAN April)
  - limited functionality and stability
- Deployment started to 10 initial sites
  - Focus not on functionality, but establishing procedures
  - Getting sites used to LCFGng
- End of August only 5 sites in
  - Lack of effort of the participating sites
  - Gross underestimation of the effort and dedication needed by the sites
    - Many complaints about complexity
    - Inexperience (and dislike) of install/config Tool
    - Lack of a one stop installation (tar, run a script and go)
    - Instructions with more than 100 words might be too complex/boring to follow
- First certified version LCG1-1\_0\_0 release September 1st (PLAN in June)
  - Limited functionality, improved reliability
  - Training paid off -> 5 sites upgraded (reinstalled) in 1 day
  - Last after 1 week....
- Security patch LCG1-1\_0\_1 first not scheduled upgrade took than 24h.
- Sites need between 3 days and several weeks to come online
  - None in not using the LCFGng setup (status Thursday)
- Now: 11 in and several Tier 2 sites waiting, 2 of the original 10 missing

**middleware was late**



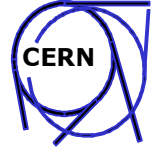
# Adding a Site



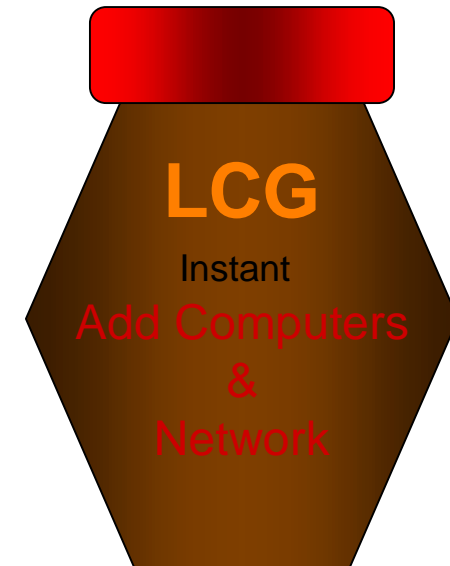
- I. Site contacts us (LCG)
  - II. Leader of the GD decides if the site can join (hours)
  - III. Site gets mail with pointers to documentation of the process
  - IV. Site fills questionnaire
  - V. We, or primary site write LCFGng config files and place them in CVS
  - VI. Site checks out config. files, studies them, corrects them, asks questions...
  - VII. Site starts installing
  - VIII. Site runs first tests locally (described in the material provided)
  - IX. Site maintains config. in CVS (helps us finding problems)
  - X. Site contacts us or primary site to be certified
    - Currently we run a few more tests, certification suite in preparation
    - Site creates a CVS tag
    - Site is added to the Information System
- We currently lack proper tool to express this in the IS



# Difficulties



- Sites without LCFGng (even using lite) have severe problems getting it right
  - We can't help too much, dependencies depend on base system installed
  - The configuration is not understood well enough (by them, by us)
  - Need one keystroke "Instant GRID" distribution (hard..)
  - Middleware's dependencies too complex
- Debugging a site
  - Can't set the site remotely in a debugging mode
  - The glue status variable covers the LRM's state
  - Jobs keep on coming
  - Discovery of the other site's setup for support is hard
- History of the components, many config files
  - No tool to pack config and send to us
  - Sites fight with FireWalls
- Some sites are in contact with grids for the 1st time
  - There is nothing like "Beginners Guide to Grids"
- LCG is on many sites not a top priority
  - Many sysadmins don't find time to work for several hours in a row
  - Instructions are not followed correctly (short cuts taken)
- Time zones slow things down a bit (The grid where the sun never sets)

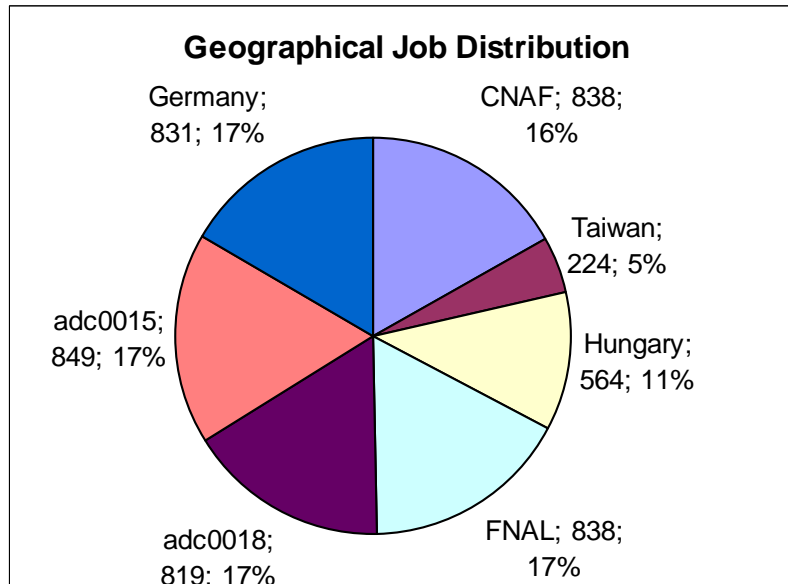




# Stability-Operation



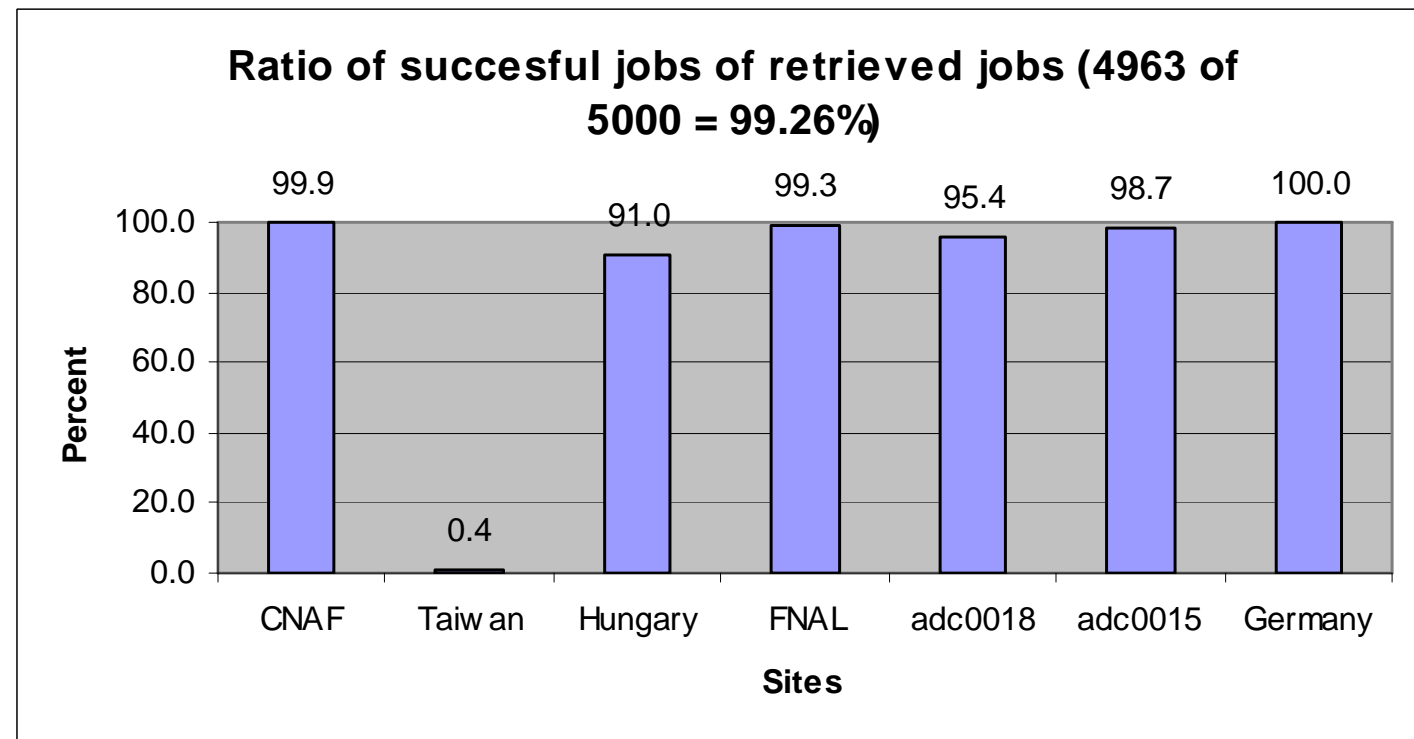
- Running jobs has now greatly improved
  - “Hello World” jobs are about 95% successful
  - Services crash with much lower rate (some bug fixes already on C&T)
  - Some bugs in LCG1-1\_0\_x already fixed on C&T
  - Grid services degrade gracefully
- So far the MDS is holding up well
- Focus in this area during the next few month
  - Long running jobs with many jobs
  - Complex jobs ( data access, many files,...)
  - Scalability test for the whole system with complex jobs
  - Chaotic (many users, asynchronous access, bursts) usage test
  - Tests of strategies to stabilize the information system under heavy load
    - We have several that we want to try as soon as more Tier2 sites join
  - We need to learn how the systems behave if operated for a long time
    - In the past some services tended to “age” or “pollute” the platforms they ran on
  - We need to learn how to capture the “state” of services to restart them on different nodes
  - Learn how to upgrade systems (RMC, LRC...) without stopping the service
    - You can’t drain LCG1 for upgrading



## LCG 1.0 Test (19./20. Sept. 2003):

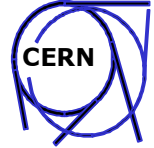
- 5 streams
- 5000 jobs in total
  - Input and OutputSandbox
  - Brokerinfo query
  - 30 sec sleep

Ingo Augustin gave this slide to me





# Next Steps



- Get everything needed for LCG2 (used in the 2004DCs) Nov. 20th
  - System for distributing experiments SW (RPMS through us won't do)
  - GCC3.2
  - VDT Globus 2.4.x
  - VOMS
    - Tests started, server is set up, will be used first for CE (Storage access later)
    - Even if we will use initially very few roles for authorization still many benefits
      - Amount of information in the IS becomes independent of number of grid users
      - Might be used to do grid wide admin. Work
      - See some problems mapping to UNIX file access (No ACL in UNIX)
  - Access to storage
    - GFAL, SRM integration with CASTOR and ENSTORE (very soon)
    - What to do with disk only sites (make them run DCASH)?
  - Distributed RLSs (can of worms -> see next slides)
  - POOL integration with distribution and RLSs
  - Basic simple accounting system (agreed on plan, found volunteers )
  - Integration of not dedicated production clusters (WNs on non routed nets)

Worms are in  
the extra s





# Next Steps II



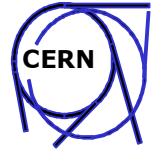
- Part2 (mainly deployment):
  - Automate installation as far as possible
  - Define installation and configuration to allow tool independent usage
  - Establish procedure to integrate Tier2 centres
    - Hierarchical model (CERN can't support 20 sites)
- Integrate US Tier2 centres
  - Middleware (based on GRID3/OSG not fully interoperable)
  - Many challenges, tests will be done together with ATLAS



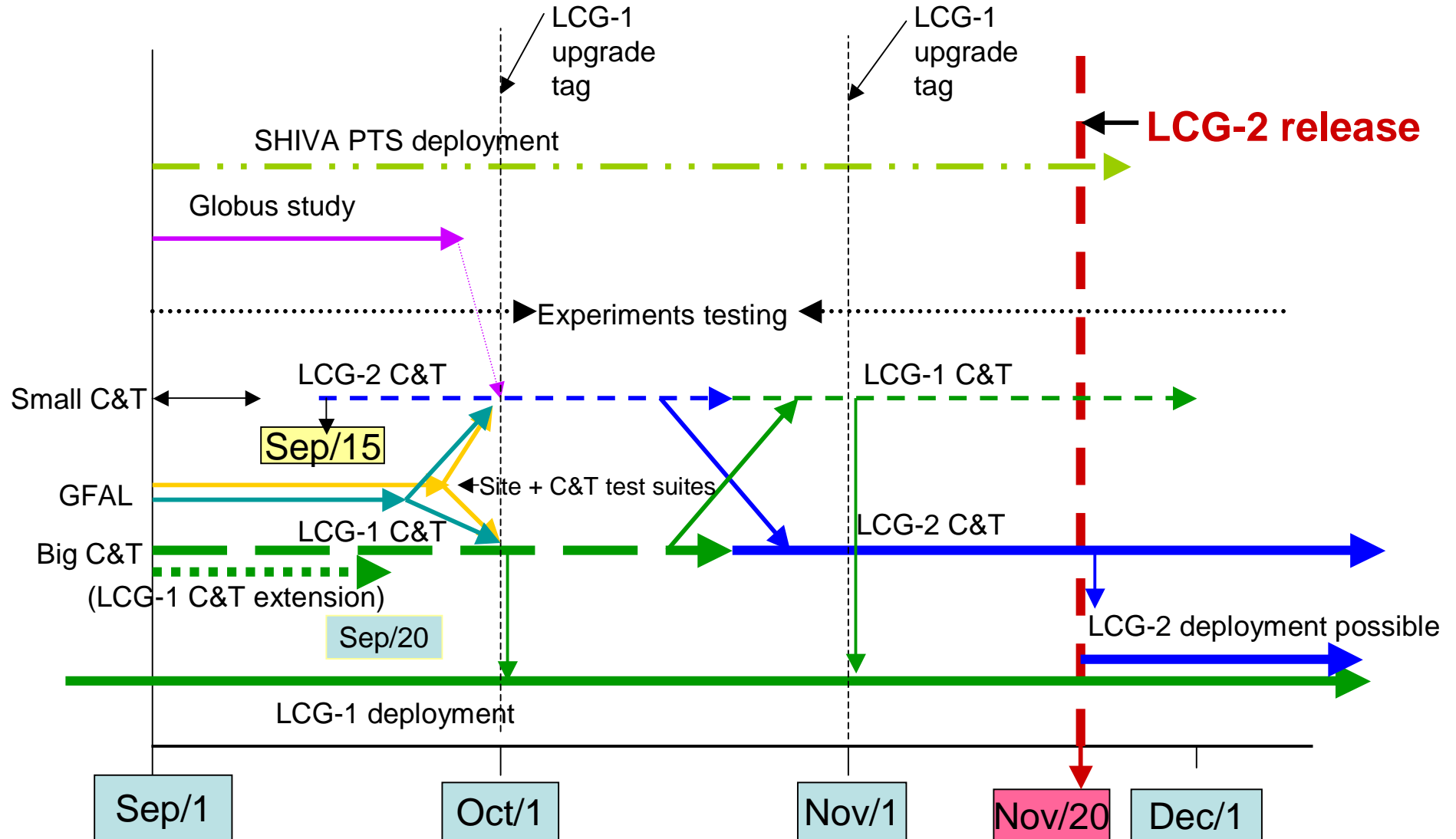
# Next III



- R-GMA
  - R-GMA still in the process of stabilization
  - Timescale:
    - Until November not testing RGMA as an replacement for MDS
    - Interoperability with US grid infrastructure is a must
    - If time permits comparison tests on the deployed system



# Timeline - Preparation





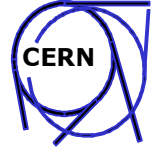
# RLSs



- Why the RLSs are a can of worms
  - In LCG two, currently non interoperable versions will be used.
  - The US Tier2 centres will deploy the Globus-RLS
  - This cuts LCG practically into two
- Plan (short version, more detailed after the summary slides)
  - Make POOL configurable to work with either of the RLSs (11/03)
  - POOL group provides tools to cross populate file catalogues (11/03)
  - Globus-RLS modified to use ORACLE as back-end
  - Integrate the RLSs (goal 5/04)
    - Integrate the Globus-RLI and EDG-LRC
    - APIs
    - Update clients (POOL, RB, ..)
- Now: Run EDG-RLS at CERN (ORACLE back-end)
  - No EDG-RLI deployed



# RLSs II



- January 2004 start with minimal solution for 2004 Data Challenges
  - Sites provide a service based either on LOCAL EDG-LRCs or Globus RLS
  - Cross-population every few hours
  - A few sites (like CERN) will get all updates and then push back out
    - Latency for the RBs has to be taken into account by production managers
    - Most production work will do anyway bulk updates at end of job
    - “Solves” requirement for a proxy replica manager. Bulk updates can be run from appropriate nodes.

This is the **current** plan.  
We have changed our plans quite often!!!



# Summary



- Middleware was 3 months late
  - Less: functionality, tests, experience with operation
  - **Core functionality**
    - Is clearly there
    - Reliability improved fundamentally compared to edg-1.4
- System now at scale foreseen (11 sites in)
  - Integration between US and European sites still an issue
- Experiments are getting ready to test the system
  - This will help to discover problems
- Very little time to turn this into a real production system
  - Critical components are just coming in (SE)
  - Has to be done incrementally on the running service
- Deploying the software at new sites not always easy
  - Different reasons (attitude, complexity, priorities, acceptance of tools)



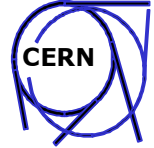
# RLS



- Issue – 2 non-interoperable implementations
- **Proposed Strategy – development**
- Integrate POOL with Globus RLS, so that it is able to communicate with both RLS implementations (but not from the same process).
  - Assuming all the above conditions are met the earliest this work could be completed would be November 2003.
  - This requires close collaboration between the POOL group and someone (a developer?) very familiar with the Globus RLS.
  - Also might require some additions/changes in the RLS API?
- POOL group provide the tools to enable cross population of POOL file catalogs between RLS implementations. These tools are basically available now.
- Globus RLS is ported to Oracle as the database back-end.
- In parallel we work on the interoperability roadmap, with a target date of May 2004 for this to be available:
  - Agree the APIs for RLS and RLI. This discussion should include agreement on the syntax of filenames in the catalog.
  - Implement the Globus RLI in the EDG RLS, make the EDG LRC talk the “Bob” protocol.
  - Implement the client APIs
  - Define and implement the proxy replica manager service
  - Update the POOL and other replica manager clients (e.g. EDG RB)



# RLS Proposed Strategy – services



- Now: Run EDG RLS at CERN and at US Tier 1 sites – at least until Globus RLS is running with Oracle.
  - The EDG RLS in this scenario is the LRC only – we will not deploy the RLI.
- By January 2004: LCG service is provided by local EDG LRCs or Globus RLS, and cross-population tools to enable catalog updates. This is the minimal solution for the 2004 Data Challenges.
  - Most batch production work will in any case use bulk updates of the catalogs rather than file-by-file updates from the job.
  - Suggest initially CERN catalog gets all updates and then pushes back out. This implies that pre-job file replication is required so that data is already at site where the job will run.
  - This model removes the requirement for a proxy replica manager since the bulk updates can be run from externally visible nodes.