# eGee

# SA1: Grid Operations and Management

*Ian Bird, CERN*

*SA1 Activity Manager*
*EGEE 2nd EU Review*
*6-7/12/2005*
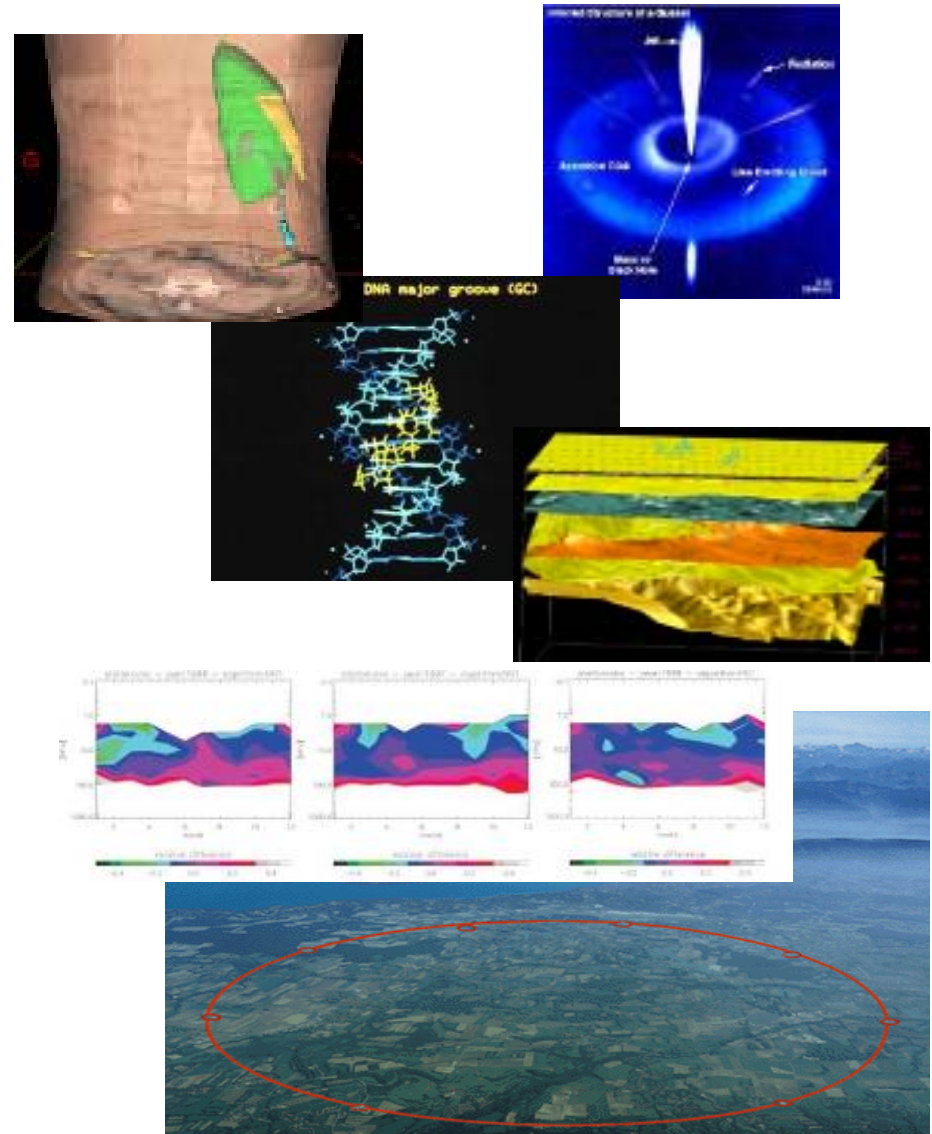
**www.eu-egee.org**

Information Society
and Media

**eGee**

Enabling Grids for E-sciencE

- **Scale and usage of infrastructure**
- **Grid Operations**
  - Metrics, operations support
- **Certification and deployment**
- **Pre-production Service**
- **User support**
- **Operational security**
- **Interoperability / interoperation**
- **Input to standards process**
- **LCG-2/gLite convergence**
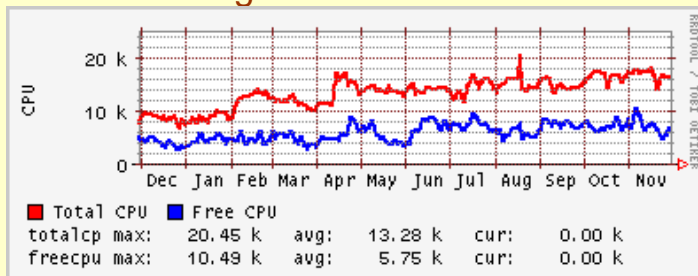- **Key points for SA1**
- **Plans for remainder of project**

➢ Many more sites than anticipated for this stage of project
- 179 actual, cf. 50 proposed for end of year 2
- ~2000 CPU in sites outside of EGEE federations (7 countries)

➢ Includes industrial partner sites (HP in Puerto Rico and UK)

➢ Exposes full complexity of grid operations – # sites not resources, nor # users

**EGEE:**

179 sites, 39 countries
>17,000 processors,
~5 PB storage



Total CPU   Free CPU
totalcp max:  20.45 k  avg:  13.28 k  cur:  0.00 k
freecpu max:  10.49 k  avg:  5.75 k   cur:  0.00 k

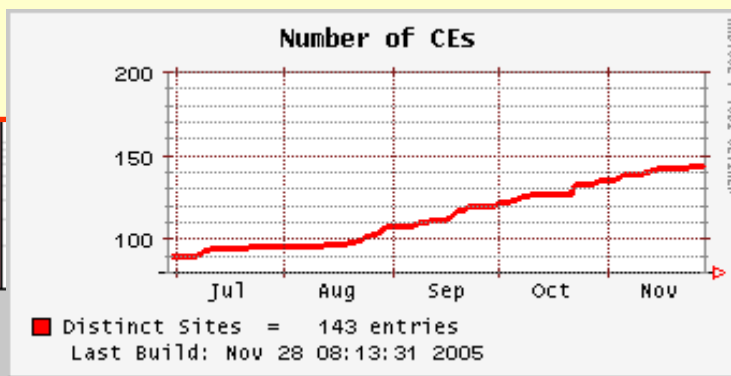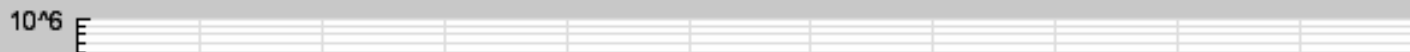| country | sites | country | sites | country | sites |
|---------|-------|---------|-------|---------|-------|
| Austria | 2 | India | 2 | Russia | 12 |
| Belgium | 3 | Ireland | 15 | Serbia | 1 |
| Bulgaria | 4 | Israel | 3 | Singapore | 1 |
| Canada | 7 | Italy | 25 | Slovakia | 4 |
| China | 3 | Japan | 1 | Slovenia | 1 |
| Croatia | 1 | Korea | 1 | Spain | 13 |
| Cyprus | 1 | Netherlands | 3 | Sweden | 4 |
| Czech Republic | 2 | Macedonia | 1 | Switzerland | 1 |
| Denmark | 1 | Pakistan | 2 | Taipei | 4 |
| France | 8 | Poland | 5 | Turkey | 1 |
| Germany | 10 | Portugal | 1 | UK | 22 |
| Greece | 6 | Puerto Rico | 1 | USA | 4 |
| Hungary | 1 | Romania | 1 | CERN | 1 |

Aggregate Accounting Plot for EGEE

□ alice
■ atlas
□ babar

From Accounting data:

→ ~3 million jobs in 2005 so far

→ Sustained daily rates (per month Jan – Nov 2005):

[2185, 2796, 7617, 10312, 11151, 9247, 9218, 11445, 10079, 11124, 9491]

→ ~8.2 M kSI2K.cpu.hours → >1000 cpu years

**Real usage is higher as accounting data was not published from all sites until recently**

Number of CEs

■ Distinct Sites = 143 entries
Last Build: Nov 28 08:13:31 2005

LCG sustained data transfers using FTS; in excess of 500 MB/s

**Domain distribution of Flexx run jobs**

**WISDOM data challenge**



**Zeus collaboration at DESY**



**ATLAS:**
**Number of jobs/day**

**June05 - Technical Design Report**

**Sep05 - SC3 Service Phase**

**May06 – SC4 Service Phase**

**Sep06 – Initial LHC Service in stable operation**

**Apr07 – LHC Service commissioned**

2005          2006          2007          2008

*SC2* ——

*SC3* ——

*SC4* ——

cosmics

First beams

First physics

Full physics run

*LHC Service Operation* ⟶

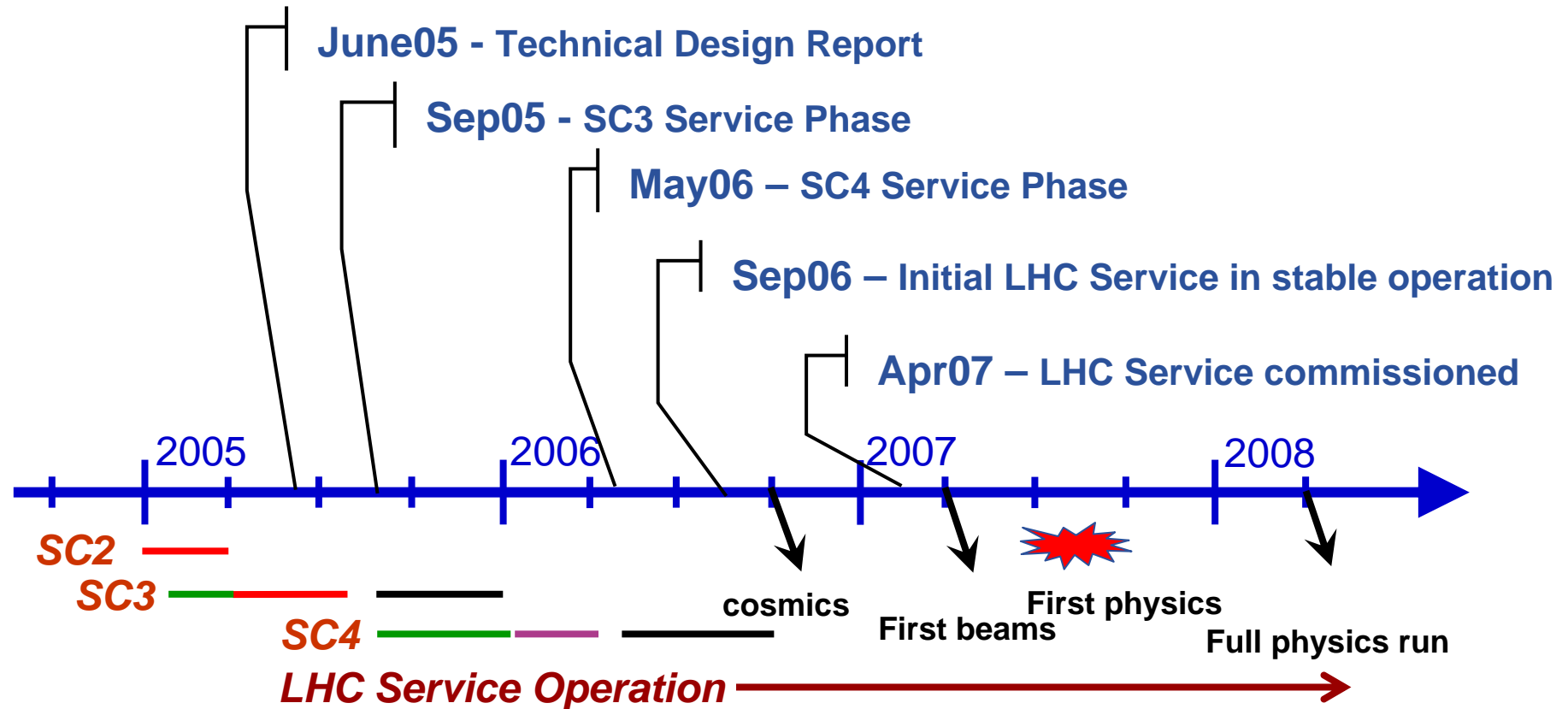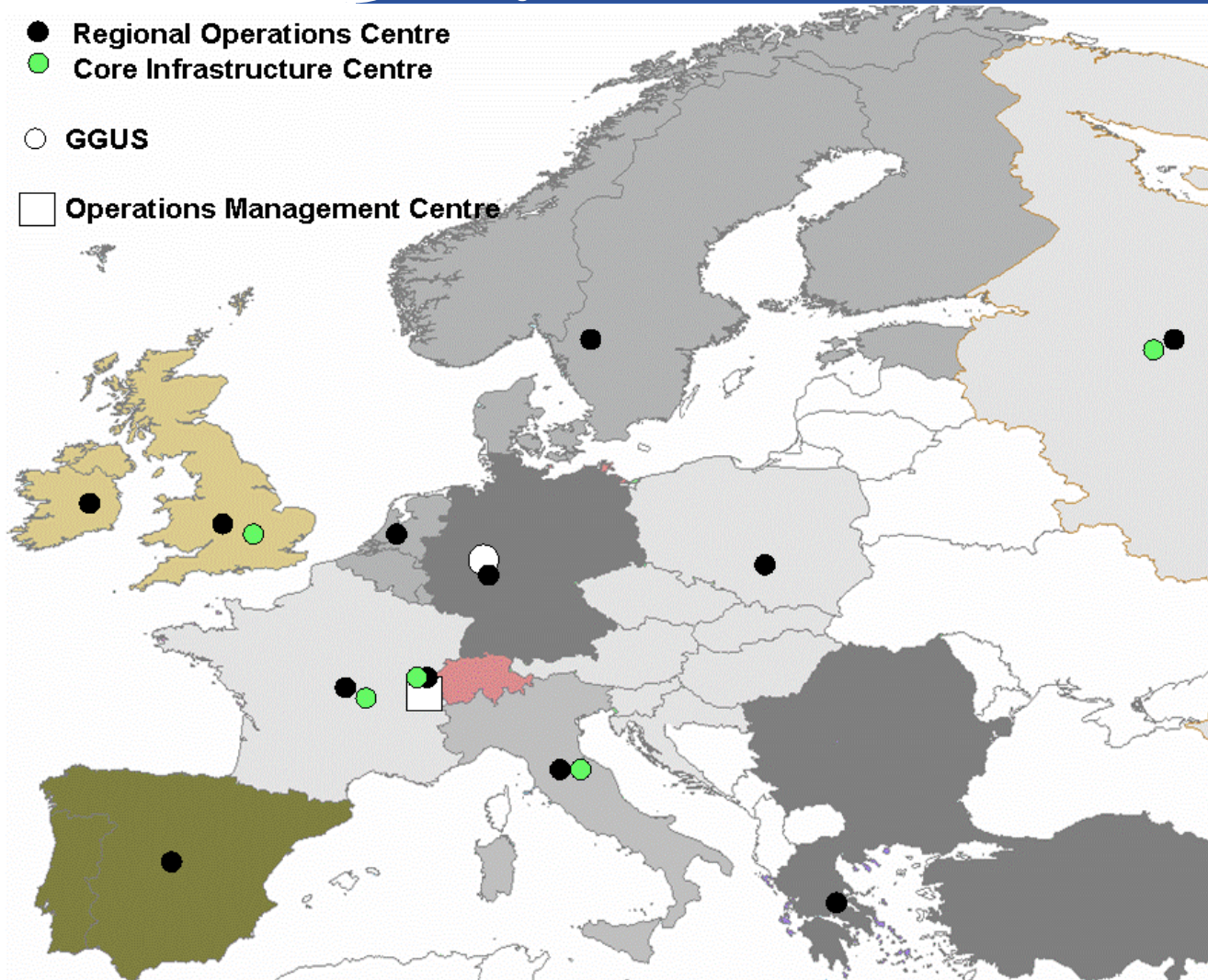**SC2** – Reliable data transfer (disk-network-disk) – 5 Tier-1s, aggregate 500 MB/sec sustained at CERN
**SC3** – Reliable base service – most Tier-1s, some Tier-2s – basic experiment software chain – grid data throughput 500 MB/sec, including mass storage (~25% of the nominal final throughput for the proton period)
**SC4** – All Tier-1s, major Tier-2s – capable of supporting full experiment software chain inc. analysis – sustain nominal final grid data throughput
**LHC Service in Operation** – September 2006 – ramp up to full operational capacity by April 2007 – capable of handling twice the nominal data throughput

- ● **Regional Operations Centre**
- ● **Core Infrastructure Centre**
- ○ **GGUS**
- □ **Operations Management Centre**

**Operations Management Centre (OMC):**

– At CERN – coordination etc

**Core Infrastructure Centres (CIC)**

– Manage daily grid operations – oversight, troubleshooting
  - ▪ "Operator on Duty"
– Run infrastructure services
– Provide 2nd level support to ROCs
– UK/I, Fr, It, CERN, Russia, Taipei

**Regional Operations Centres (ROC)**

– Front-line support for user and operations issues
– Provide local knowledge and adaptations
– One in each region – many distributed

**User Support Centre (GGUS)**

– In FZK: provide single point of contact (service desk), portal

**eGee**

**Enabling Grids for E-sciencE**

- **CIC – on – duty (grid operator on duty)**
  - Started November 2004
  - 6 teams working in weekly rotation
    - CERN, IN2P3, INFN, UK/I, Ru,Taipei
  - Crucial in improving site stability and management
- **Operations coordination**
  - Weekly operations meetings
  - Regular ROC, CIC managers meetings
  - Series of EGEE Operations Workshops
    - Nov 04, May 05, Sep 05
    - Last one was a joint workshop with Open Science Grid
  - These have been extremely useful
    - Will continue in Phase II
    - Bring in related infrastructure projects – coordination point
    - Continue to arrange joint workshops with OSG (and others?)
- **Geographically distributed responsibility for operations:**
  - There is no "central" operation
  - Tools are developed/hosted at different sites:
    - GOC DB (RAL), SFT (CERN), GStat (Taipei), CIC Portal (Lyon)
- **Procedures described in Operations Manual**

- **Improvement in site stability and reliability is due to:**
  - **CIC on duty oversight and strong follow-up**
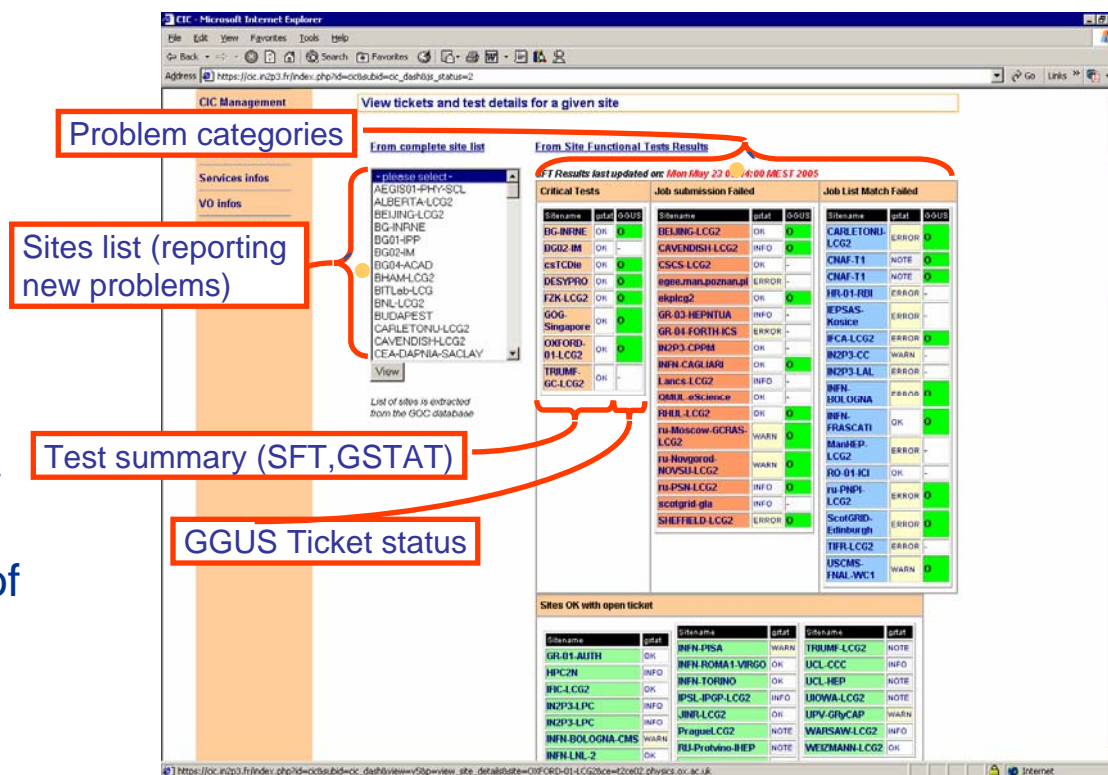  - **Site Functional Tests, Information System monitor**

- **Many complementary monitoring tools, 2 important tools:**
  - Site Functional Tests (SFT)
  - Information System monitor (GStat)
- **Dashboard provides top level view of problems:**
  - Integrated view of monitoring tools (summary) - shows only failures and assigned tickets
  - Detailed site view with table of open tickets and links to monitoring results
  - Single tool for ticket creation and notification emails with detailed problem categorisation and templates
  - Ticket browser with highlighting expired tickets



- **Well maintained – is adapted quickly to new requirements/suggestions**

**eGee**

Enabling Grids for E-sciencE

- **Site Functional Tests (SFT)**
  - Framework to test services at all sites
  - Shows results matrix
  - Detailed test log available for troubleshooting and debugging
  - History of individual tests is kept
  - Can include VO-specific tests (e.g. sw environment)
  - SFT's have evolved to become stricter as lessons are learned
  - Normally >80% of sites pass SFTs
    - NB of 180 sites, some are not well managed
- **Freedom of Choice tool (FCR)**
  - Uses results of SFT
  - Allows apps to select good sites according to their criteria
  - Selection of "critical" tests for each VO to define which sites are good/bad
  - VO can select set of functional tests that it requires
  - Can white- or black-list sites
  - Operator can remove site from production
- **SFT framework and FCR tool provide dynamic selection of "good" sites**

- Very important in stabilising sites:
  - Apps use only good sites
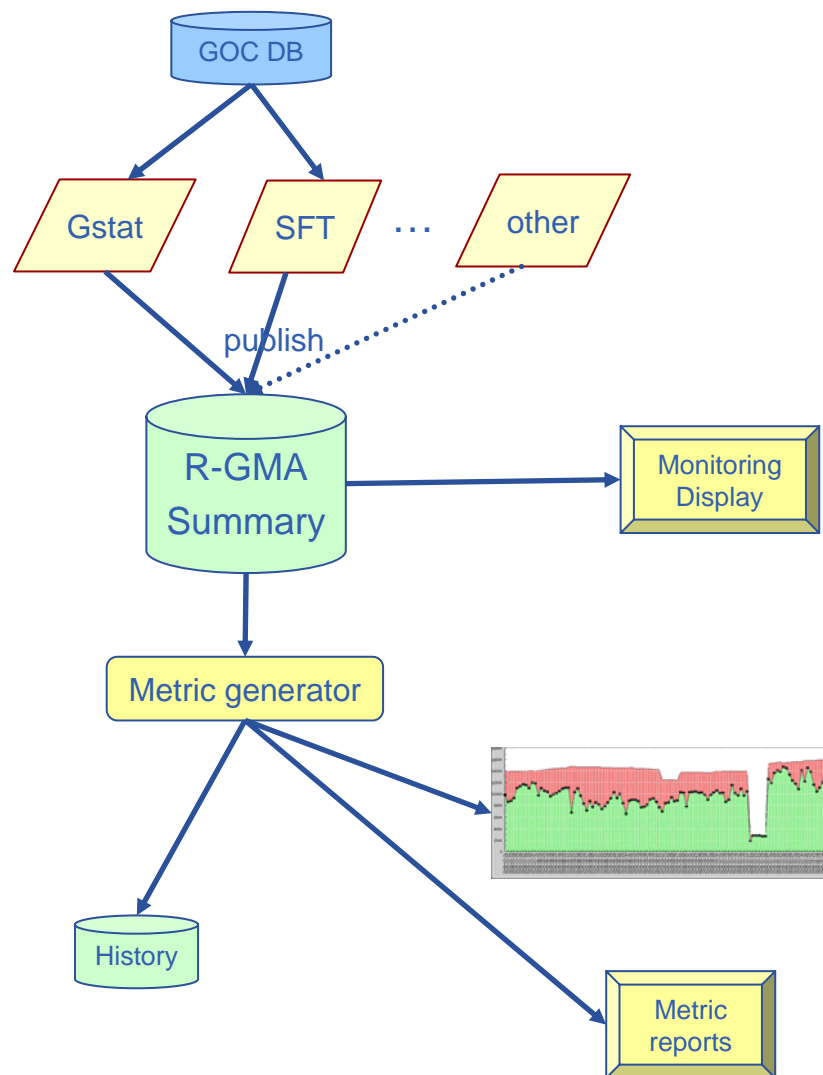  - Bad sites are automatically excluded
  - Sites work hard to fix problems

| | crl | CRL timestamp test |
|---|---|---|
| JS  Job submission failed  #f4876b | rm | Replica Management |
| CT  Critical tests failed  #f9d48e | votag | VO Tag management |
| NT  Non-critical tests failed  #f2f98e | js | Job submission |
| OK  OK  #b2f98e | bi | BrokerInfo |

**Test summary**

| | SD | JL | JS | CT | OK | total |
|---|---|---|---|---|---|---|
| dteam | 15 | 12 | 4 | 6 | 139 | 176 |
| lhcb | 15 | 81 | 5 | 35 | 39 | 175 |

| | St. | Site Name | Site CE | VO dteam St. | js | ver | wn | ca | rgma | bi | csh | rm | votag | swdir | crl | VO lhcb St. | js | dirac-test |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | **AsiaPacific** | | | | | | | | | | | | | | | | |
| 1. | CT | INDIACMS-TIFR | ce.indiacms.res.in | CT | O | 2 6 0 | I | O | O | O | X | O | | O | !!! | JL | X | ?? |
| 2. | OK | TW-NCUHEP | grid01.phy.ncu.edu.tw | OK | O | 2 6 0 | I | O | O | O | O | O | | O | !!! | JL | X | ?? |
| 3. | OK | TOKYO-LCG2 | dgce0.icepp.jp | OK | O | 2 4 0 | I | O | O | O | O | O | | O | !!! | JL | X | ?? |
| 4. | OK | Taiwan-LCG2 | lcg00125.grid.sinica.edu.tw | OK | O | 2 6 0 | I | O | O | O | O | O | | O | !!! | JL | X | ?? |
| 5. | OK | Taiwan-IPAS-LCG2 | testbed001.phys.sinica.edu.tw | OK | O | 2 6 0 | I | O | O | O | O | O | | O | !!! | JL | X | ?? |
| 6. | OK | GOG-Singapore | melon.ngpp.ngp.org.sg | OK | O | 2 6 0 | I | O | O | O | O | O | | O | !!! | JL | X | ?? |
| 7. | OK | Taiwan-NCUCC-LCG2 | ce.cc.ncu.edu.tw | OK | O | 2 6 0 | I | O | O | O | O | O | | O | !!! | OK | O | O |
| 8. | OK | LCG_KNU | cluster50.knu.ac.kr | OK | O | 2 5 0 | I | O | O | O | O | O | | O | !!! | CT | O | !!! |
| | | **BNL** | | | | | | | | | | | | | | | | |
| 9. | SD | BNL-LCG2 | lcg-ce01.usatlas.bnl.gov | SD | X | ?? | ?? | ? ? | ?? | ? ? | ?? | ?? | ?? | ?? | SD | X | ?? |
| | | **Canada** | | | | | | | | | | | | | | | | |
| 10. | JL | TORONTO-LCG2 | bigmac-lcg-ce.physics.utoronto.ca | JL | X | 2 6 0 | I | O | O | O | O | O | | W | O | !!! | OK | O | O |
| 11. | SD | CARLETONU-LCG2 | lcg02.physics.carleton.ca | SD | X | ?? | ?? | ? ? | ?? | ? ? | ?? | ?? | ?? | ?? | SD | X | ?? |
| 12. | OK | TRIUMF-LCG2 | lcgce01.triumf.ca | OK | O | 2 6 0 | I | O | O | O | O | O | | O | O | OK | O | O |
| 13. | OK | Umontreal-LCG2 | lcg-ce.lps.umontreal.ca | OK | O | 2 6 0 | I | O | O | O | O | O | | W | O | !!! | OK | O | O |

**eGee**

Enabling Grids for E-sciencE

- **R-GMA is used as the "universal bus" for monitoring information**
- **SFT and GStat both publish results to R-GMA using common schema**
- **GOC DB source of:**
  - Sites and nodes to monitor,
  - Status (downtime, etc.)
- **Scalability:**
  - Currently >170 sites
  - About 3.5M tuples for 1 month history with full detail
  - After one month only summary information
- **Aggregate views →**
  - Dashboard, high level monitors
  - Eventually automated alarms
- **Summary information**
  - Generate metrics: site availability
- **Framework – longer term**
  - Include results from various tools
  - Aggregate the disparate data
  - Generate alarms

Available sites weekly (2005-11-11)

Available cpus weekly (2005-11-11)

Available sites

Daily: July → November

| Service | Class | Comment |
|---------|-------|---------|
| SRM 2.1 | C | Monitoring of SE |
| LFC | C/H | |
| FTS | C | Base on SC experience |
| CE | C | Monitored by SFT now |
| RB | C | Job monitor exists |
| Top level BDII | C | Can be included in Gstat |
| Site BDII | H | Monitored by Gstat |
| MyProxy | C | |
| VOMS | H | |
| R-GMA | H | |

C: Critical service

H: High availability

**Effort identified in various ROCs to provide availability tests for each service**

**Will all be integrated into SFT framework**
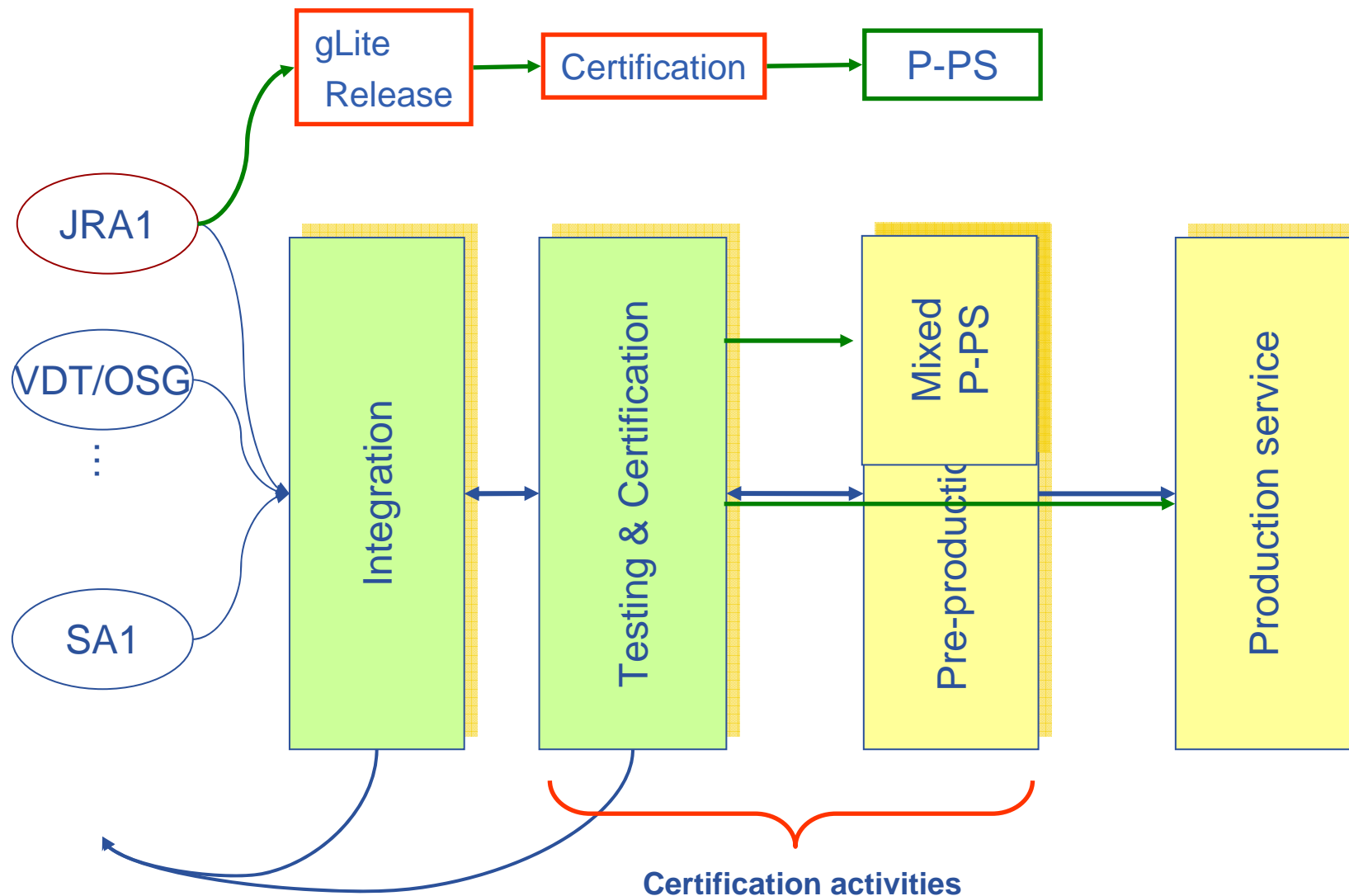
**First approach to SLA:**
- **each Class (C, H, etc) defines required service availability**

# Checklist for a new service

**egee**
Enabling Grids for E-sciencE

- **User support procedures (GGUS)**
  - Troubleshooting guides + FAQs
  - User guides
- **Operations Team Training**
  - Site admins
  - CIC personnel
  - GGUS personnel
- **Mon**
- **Acc**
- **Service Parameters**
  - Scope - Global/Local/Regional
  - SLAs
  - Impact of service outage
  - Security implications
- **Contact Info**
  - Developers
  - Support Contact
  - Escalation procedure to developers
- **Interoperation**
  - Documented issues

- **First level support procedures**
  - How to start/stop/restart service
  - How to check it's up
  - Which logs are useful to send to CIC/Developers
    - and where they are

> ➤ **This is what is takes to make a reliable production service from a middleware component**

> ➤ **Not much middleware is delivered with all this ... yet**

- **Tools for CIC to spot problems**
  - GIIS monitor validation rules (e.g. only one "global" component)
  - Definition of normal behaviour
    - Metrics
- **CIC Dashboard**
  - Alarms
- **Deployment Info**
  - RPM list
  - Configuration details
  - Security audit

- **Deployment process has improved significantly:**
  - Significant effort to improve the deployment process – better separation of functional improvements from critical updates
  - Simplified installation and configuration tools – made life much simpler for administrators
  - Wider deployment testing before release; also pre-production
  - GGUS coordinates problem follow up
- **Certification:**
  - Increased effort was identified (UK, INFN) to address lack of testing of new gLite components
  - Parallel processes to speed up gLite testing:
    - Production certification
    - "pure" gLite certification
    - Mixed (LCG-2.x + gLite) → this will become primary strategy
  - gLite 1.4.1 is being certified now

gLite Release → Certification → P-PS

JRA1

VDT/OSG

...

SA1

Integration

Testing & Certification

Mixed P-PS

Pre-production

Production service

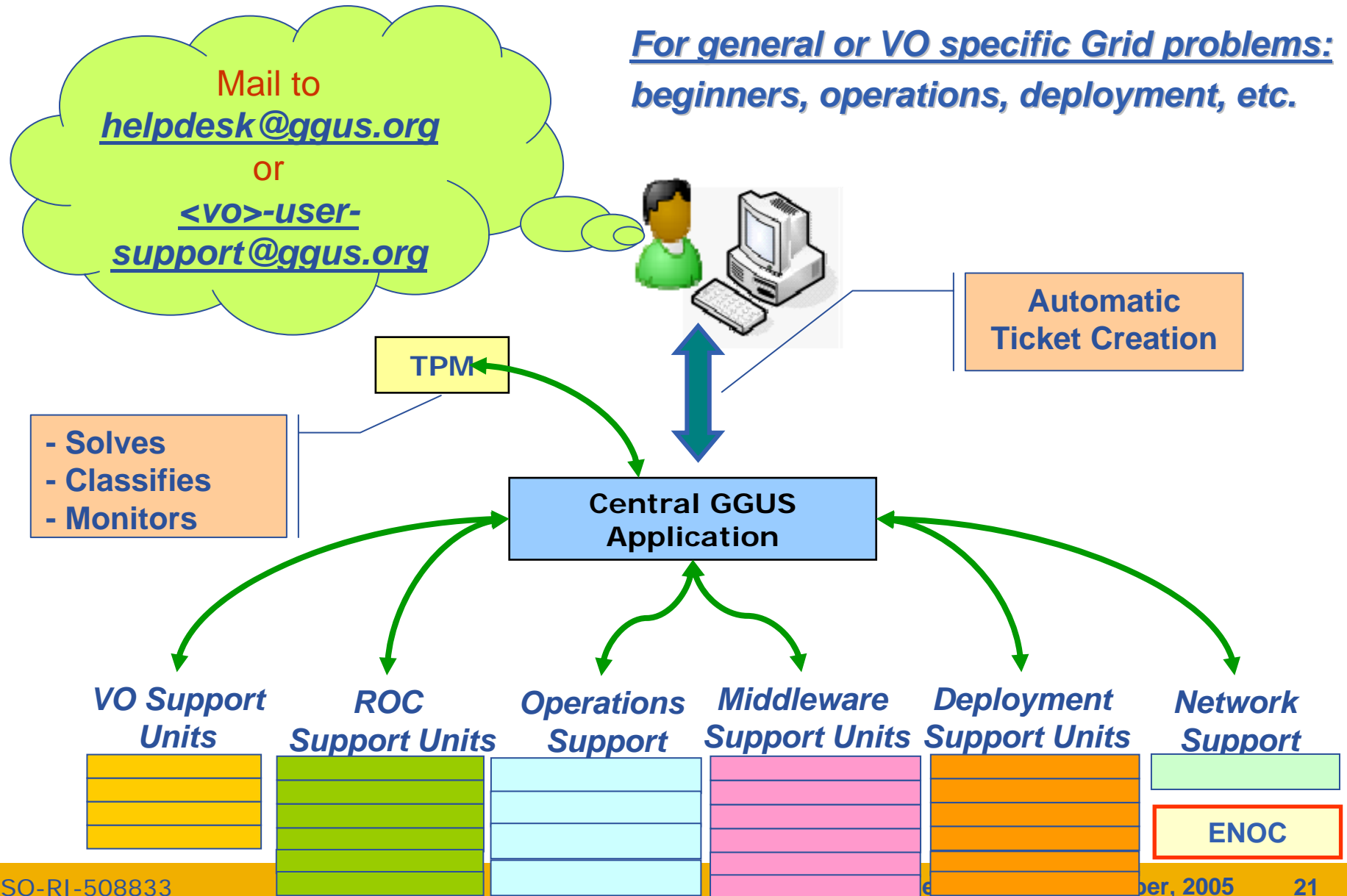**Certification activities**

**eGee**

Enabling Grids for E-sciencE

- **Goals: (tension between …)**
  - Applications: Rapid updates – new functionality, bug fixes
  - Sites: Fixed schedules – to ensure good planning and response
  - Deployment team: Sufficient certification and testing – to ensure quality
  - Rapid reaction by sites to new releases is desired by applications and deployment team
- **Lessons learned**
  - **EGEE production service is a grid of independent federations**
    - **ROCs schedule upgrades in their region**
    - New releases need a few months to reach 80% site deployment
  - Early announcement of new releases needed
    - To allow time for external deployment testing ($\rightarrow$ p-ps)
  - Release definition non-trivial with 3 months intervals
    - Closing door for changes is almost impossible
  - Certification Tests need to be extended (performance tests)
  - Patches have to come with a standard set of information
    - Ports, variables, config changes …
  - Updates work quite well
- **Now: Integrate JRA1 and SA1 processes**
  - Take into account the experiences gained over past 4 years
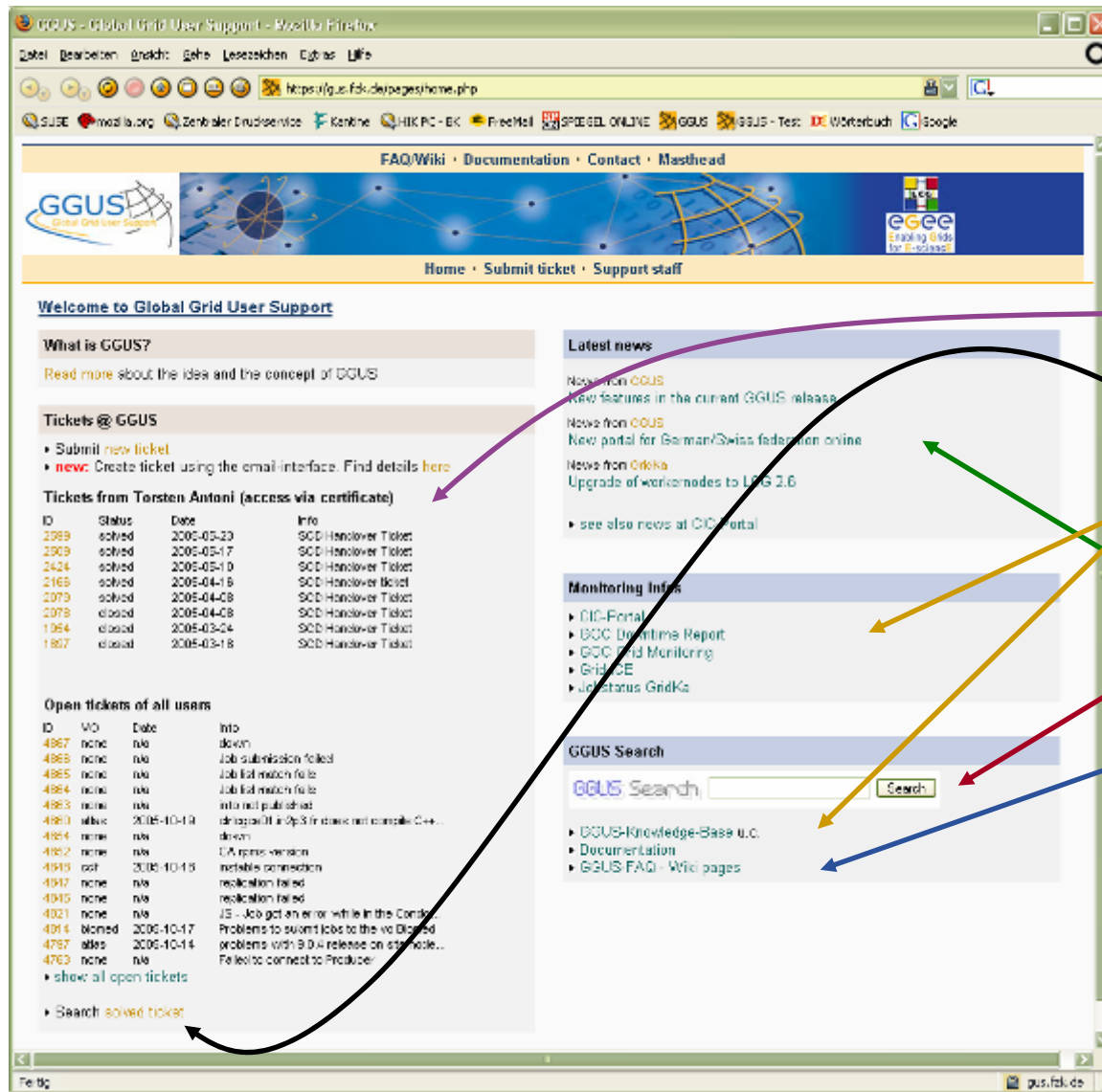  - Ensure (TCG) priorities are driven by the applications

**eGee**

Enabling Grids for E-sciencE

- **Current P-PS is a "pure" gLite service**
  - BDII, SRM SE and MyProxy server are also needed.

- **The P-PS is available and used by many VOs**
  - HEP VOs (CMS, ATLAS, Alice, LHCb)
  - ARDA
  - BioMed
  - egeode
  - NA4 (testing)
  - DILIGENT
  - SWITCH

- **Currently upgrading from gLite 1.4 to gLite 1.4.1 (a major patch)**
  - As the service is now in use, upgrades are planned and phased to minimize the impact to users.

- **Currently preparing to move the day-to-day operations of the P-PS to the production operations team**
  - SFT monitoring is now in place
  - All P-PS sites are now correctly entered in the GOC database
  - Production operation processes are being implemented for the P-PS

- **Planning is under way for moving the P-PS from being a pure gLite service to being a true pre-production service which closely mirrors production**

**eGee**

**Enabling Grids for E-sciencE**

| ROC | Site | CPUs | SE | Core Services | | | |
|-----|------|------|-----|------|------|------|------|
| Asia-Pacific | ASGC | 0 | | WMS | | | |
| CE | CYFRONET | 3 | | | | | |
| CERN | CERN | 54 | DPM | WMS | FTS | VOMS (production) | |
| DE/CH | FZK | 2 | | | | | |
| France | IN2P3 | 4 | | | FTS | VOMS | |
| Italy | CNAF | 150 | DPM | WMS | | VOMS | BDII |
| Italy | INFN-Padova | 100 | | | | | |
| NE | NIKHEF | 0 | | | | VOMS | |
| SEE | UoM | 2 | | | | | |
| SEE | UPATRAS | 3 | | WMS | | | |
| SWE | CESGA | 2 | | | | | R-GMA |
| SWE | IFIC | 1 | Castor | | | | |
| SWE | LIP | 2 | DPM | | | | MyProxy |
| SWE | PIC | 180 | Castor | WMS | | | FireMan |
| UK/I | ScotGrid-Glasgow | 0 | | | FTS | | |

- **PIC, CNAF, Padova and CERN have given access to production batch farms**
  - PIC, Padova and CNAF running LCG WNs; CERN running gLite WNs.
  - CERN: queue to production batch farm is currently restricted to 50 jobs. This restriction can be removed, increasing the number of CPUs at CERN to ~1,500.
- **To date, over 1.5 million jobs have been submitted to the P-PS WMSs.**

**Enabling Grids for E-sciencE**

- **User Support in EGEE (helpdesk, call-centre)**
  - Regional support with central coordination (GGUS @ FZK)
  - GGUS platform connects:
    - CICs, ROCs, VOs, service teams providing support
    - Middleware developers and support
    - Networking activities (training etc).
  - Ticket Process Managers – oversee problem lifecycle
    - Ensure problems assigned and followed up
    - Problem resolution by volunteer experts – harness informal processes
  - Users can report via local helpdesks, ROC helpdesk, VO helpdesk, or to GGUS
  - Ticket traffic increasing
    - Now: Change in users from a few, experienced, production managers to general users (low quality of tickets)
- **VO support**
  - Other aspect of user support – direct support to apps to integrate with grid middleware
  - Application driven process: set up several task forces to implement this (follow successful model in LCG)
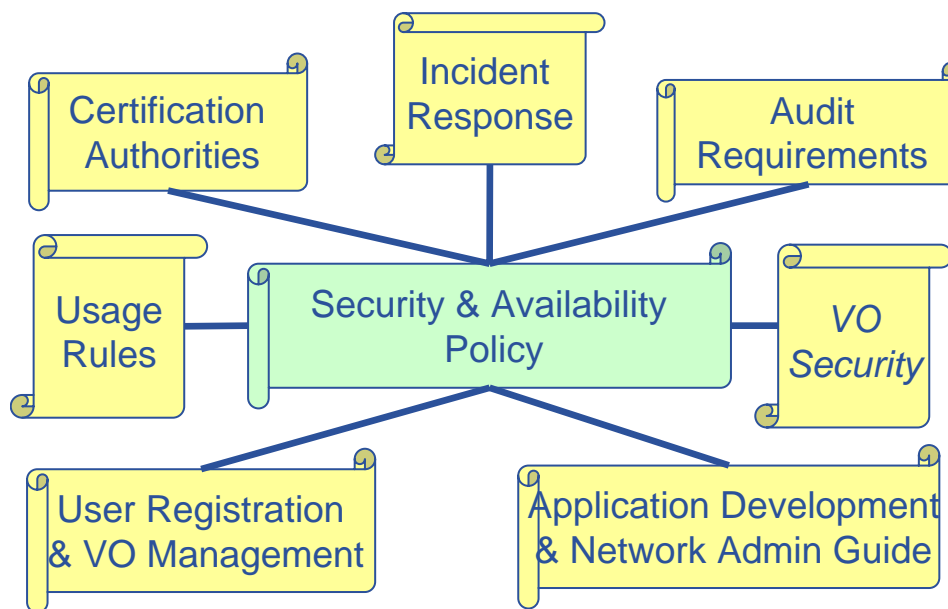
**eGee**

**Enabling Grids for E-sciencE**

*For general or VO specific Grid problems:*
*beginners, operations, deployment, etc.*

Mail to
*helpdesk@ggus.org*
or
*<vo>-user-support@ggus.org*

**Automatic Ticket Creation**

**TPM**

**- Solves**
**- Classifies**
**- Monitors**

**Central GGUS Application**

*VO Support Units*

*ROC Support Units*

*Operations Support*

*Middleware Support Units*

*Deployment Support Units*

*Network Support*

**ENOC**

**Browseable tickets**

**Search through solved tickets**

**Useful links (Wiki FAQ)**

**Latest News**

**GGUS Search Engine**

**Updated documentation (Wiki FAQ)**

- **Joint Security Policy Group**
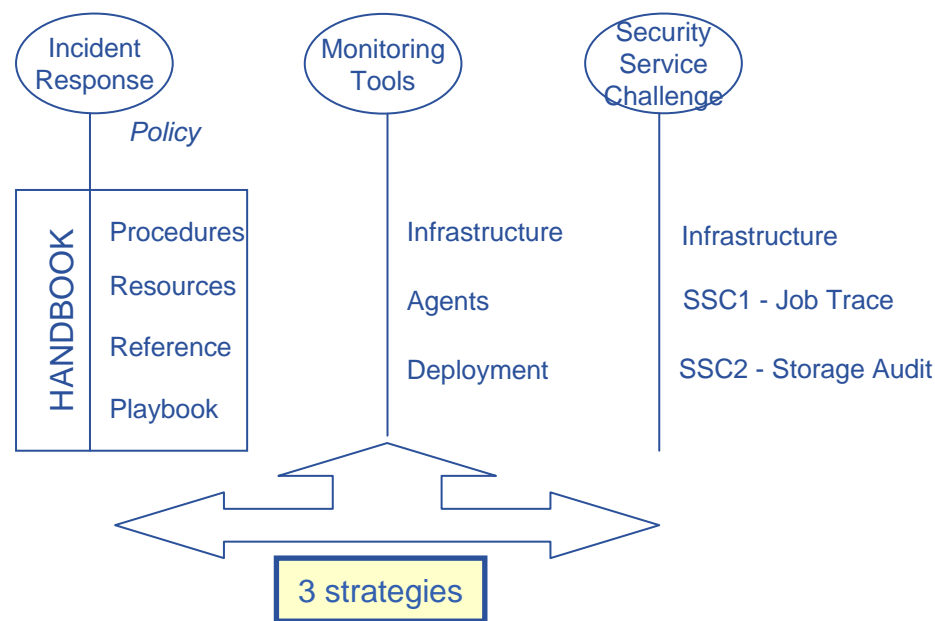  - EGEE with strong input from OSG
  - Policy Set:



- **Policy Revision In Progress/Completed**
  - Grid Acceptable Use Policy (AUP)
    - https://edms.cern.ch/document/428036/
    - common, general and simple AUP
    - for all VO members using many Grid infrastructures
      - *EGEE, OSG, SEE-GRID, DEISA, national Grids…*
  - VO Security
    - https://edms.cern.ch/document/573348/
    - responsibilities for VO managers and members
    - VO AUP to tie members to Grid AUP accepted at registration
  - Incident Handling and Response
    - https://edms.cern.ch/document/428035/
    - defines basic communications paths
    - defines requirements (MUSTs) for IR
      - *reporting*
      - *response*
      - *protection of data*
      - *analysis*
    - not to replace or interfere with local response plans

- **OSCT membership → ROC security contacts**
  - What it is not:
    - Not focused on middleware security architecture
    - Not focused on vulnerabilities (see *Vulnerabilities Group*)
  - Focus on Incident Response Coordination
    - Assume it's broken, how do we respond?
    - Planning and Tracking
  - Focus on 'Best Practice'
    - Advice
    - Monitoring
    - Analysis
  - Coordinators for each EGEE ROC
    - plus OSG LCG Tier 1 + Taipei



Incident Response — *Policy* — HANDBOOK: Procedures, Resources, Reference, Playbook

Monitoring Tools — Infrastructure, Agents, Deployment

Security Service Challenge — Infrastructure, SSC1 - Job Trace, SSC2 - Storage Audit

3 strategies

- **Has been set up this summer (CCLRC lead)**

- **Purpose: inform developers, operations, site managers of vulnerabilities as they are identified and encourage them to produce fixes or to reduce their impact**

- **Set up (private!) database of vulnerabilities**

  - To inform sites and developers

- **Urgent action → OSCT to manage**

- **After reaction time (45 days)**

  - Vulnerability and risk analysis given to OSCT to define action – publication?

  - Will not publish vulnerabilities with no solution

- **Intend to report progress and statistics on vulnerabilities by middleware component and response of developers**

- **Balance between open responsible public disclosure and creating security issues with precipitous publication**

**Enabling Grids for E-sciencE**

- **EGEE – OSG:**
  - Job submission demonstrated in both directions
  - Done in a sustainable manner
  - EGEE BDII and GIP deployed at OSG sites
    - Will also go into VDT
  - EGEE client tools installed by a grid job on OSG sites
    - Small fixes to job managers to set up environment correctly
- **EGEE – ARC:**
  - 2 workshops held (September, November) to agree strategy and tasks
  - Longer term: want to agree standard interfaces to grid services
  - Short term:
    - EGEE→ARC: Try to use Condor component that talks to ARC CE
    - ARC→EGEE: discussions with EGEE WMS developers to understand where to interface
  - Default solution: NDGF acts as a gateway

**Enabling Grids for E-sciencE**

- **Goal: to improve level of "round-the-clock" operational coverage**

- **OSG have been to all of the EGEE operations workshops**
  - Latest was arranged as a joint workshop

- **Can we share operational oversight?**
  - Gain more coverage (2 shifts/day)

- **Share monitoring tools and experience**
  - Site Functional tests (SFT)
  - Common application environment tests
  - Work on common schema for monitoring data started

- **User support workflows – interface**

- **Strong interest from both sides**


- **Now: Write a short proposal of what we can do together**
  - Both EGEE and OSG have effort to work on this

- **Follow up in future operations workshops**

**Enabling Grids for E-sciencE**

- **Interoperation and interoperability**
  - De-facto standards – common understandings/interfaces
    - GT2, GSI, SRM, BDII/GIP (MDS), …
  - Agreement on schema:
    - GLUE 1.2/GLUE 2.0; GGF Usage record for accounting
      - *GLUE 2.0 will unify EGEE, OSG, ARC information schema*
    - Consider: common operations and job monitoring schema
- **Top-down vs bottom-up standards – must keep a balance in production**
  - What is working now (SRM, GLUE) vs what will help in future
  - Must maintain production service while introducing new components that apply standards → slow
- **Operations:**
  - SA1 "Cookbook": summary of choices and experience deploying EGEE → intend to publish to GGF production grids
  - All aspects of operational security are very much collaborative with OSG and others (and very active in GGF)
  - Integration and certification is hard – standard interfaces and protocols should help
- **Operations Workshops**
  - Open to related infrastructure projects (EELA, EUMedGrid, SEE-Grid, … OSG, etc.)
  - Provide practical standardisation forum for which no equivalent in GGF as yet
- **SC05 Interoperability discussions**
  - Integrate bi-lateral interoperability work
  - EGEE/SA1 will contribute its work and experiences

**eGee**

Enabling Grids for E-sciencE

- **The current production middleware ("LCG-2") is stable and is daily heavily used**
  - This has to be maintained as new components are added or components replaced
  - This will always be the case – there will always be new or better services coming
  - Thus, the production distribution must evolve in a controlled way that does not break existing applications but that adds new, or improves existing, functionality
- **There is a strong and reliable process in place**
  - Integration, testing, certification, pre-production, production
  - Process constantly evaluated and improved
  - All significant components of gLite 1.4 are either in production (R-GMA, VOMS, FTS) …
  - … or on the pre-production service (CE, WMS, Fireman, gliteIO)
  - Anticipate these being available in production distributions (alongside existing components at first) – by mid-2006 (many sooner)
- **The current LCG and gLite middleware will converge to a single _distribution_ called gLite in early 2006**
- **Should not expect (or desire!) a big-bang switch to gLite (or anything else)**
- **Deploying in production any new software is a slow and time-consuming process, this lesson has been learned many times**

- **Accomplishments:**
  - SA1 is operating world's largest grid infrastructure for science
  - Significant resources available
  - In use by many real production applications
    - 10K jobs/day
  - Daily operations model is now well established
  - User support process is in place and being used
    - But it is complex !
  - Site stability is better controlled
    - Apps can select good sites
    - Understanding of metrics and what SLA might look like
  - Ports to other architectures now exist
    - IA64, other Linuxes
  - Convergence of middleware stacks under way
    - gLite components reaching production

- **Issues:**
  - Hard to balance:
    - Needs of applications for rapid updates
    - Reliable scheduling wanted by sites
    - Adequate testing and certification
  - Moving new middleware into production is time consuming:
    - Unrealistic expectations
    - Very stressful
    - But software industry knows …
  - Essential to maintain stable production environment
    - While introducing new functionality, new services
    - Backwards compatibility
    - Expensive in resources and support
  - Release of accounting (& other) data
    - some site policies restrict release of per-user data (privacy laws)
    - Accounting, job monitoring, …
  - Introducing new VOs is still too difficult

**Enabling Grids for E-sciencE**

- **Remainder of EGEE**
  - Milestones:
    - MSA1.5 (PM21) – Expanded production grid available (50 sites)
  - Deliverables:
    - DSA1.7 (PM19) – Cookbook – internal review
    - DSA1.8 (PM23) – Assessment of production operation (update of DSA1.4)
    - DSA1.9 (PM21) – Release notes corresponding to MSA1.5
  - Full metrics programme implemented (scope agreed in Pisa, Oct '05)
    - Service availability SLA for LCG (MoU)
  - Deploy major gLite components in production
- **Sustainability**
  - Prepare processes for EGEE-II
  - Re-focus on middleware support and building deployable distributions: Merge integration, testing (JRA1) with integration and certification (SA1) into single team with distributed partners
  - Work with embryonic TCG to ensure application driven priorities reflected in development and deployment priorities

**eGee**

Enabling Grids for E-sciencE

- **Infrastructure at a scale much larger than anticipated for end of year 2:**
  - 179 sites, 17k CPU, 39 countries
- **Being used at a significant scale for daily production:**
  - Sustaining > 10k jobs per day over many months
  - Many applications, not just HEP
  - Massive sustained data throughput > 500MB/s for 10 days
  - LCG service challenges, Biomed (WISDOM) data challenge
- **Operational oversight – grid "operator on duty"**
  - In place for 1 year, CERN, IN2P3, INFN, CCLRC, Russia, ASGC
  - Improved stability of sites → VO-specific selection of "good" sites
  - Metrics on stability and availability → SLAs
- **Pre-production service available**
  - In use by many applications, as testing ground for gLite
- **gLite components now in deployed middleware distribution**
  - VOMS, R-GMA, FTS, others (WMS) being certified now
- **Interoperability**
  - With OSG demonstrated, work in progress with ARC
  - Shared operational oversight with OSG under discussion