



# CCRC'08

*Jeff Templon*

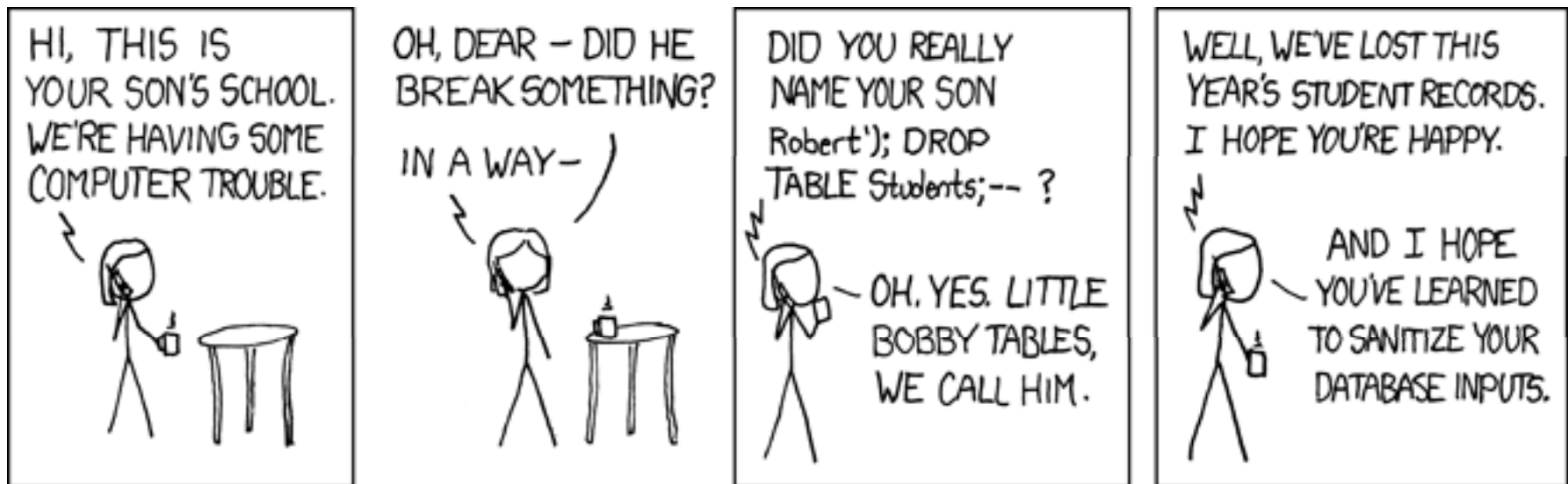
*NIKHEF*

JRA1 All-Hands Meeting

Amsterdam, 20 feb 2008



# At least our users aren't malicious



# What happens when

- ◆ Each experiment streams data into the T0
- ◆ The experiments' data model is followed T0-T1-T2
- ◆ Necessary computing (reconstruction, calibration) is done
- ◆ Sites try to reach the MoU targets for uptimes
- ◆ Graduate students try to analyze the data as it comes in
- ◆ All four experiments try this at the same time, at scale

CCRC

# CCRC Goals

- ◆ For February, we have the following three sets of metrics:
  1. The **scaling factors** published by the experiments for the various **functional blocks** that will be tested. These are monitored continuously by the experiments and reported on at least weekly;
  2. The lists of **Critical Services**, also defined by the experiments. These are complementary to the above and provide additional detail as well as service targets. It is a goal that all such services are handled in a standard fashion - i.e. as for other IT-supported services - with appropriate monitoring, procedures, alarms and so forth. Whilst there is no commitment to the problem-resolution targets - as short as 30 minutes in some cases - the follow-up on these services will be through the daily and weekly operations meetings;
  3. The services that a site must offer and the corresponding availability targets based on the **WLCG MoU**. These will also be tracked by the operations meetings.
- Phase 2 of CCRC in May

# Scaling Factors

- ◆ We won't make it in Feb ... Functional problems
  - Castor can't handle load
  - At least on exp't framework can't handle load
  - FTS corrupt proxy problems (race condition)
- ◆ Note how "data driven" HEP is : without functional data flow, test at scale is not possible!

# FTS “corrupted proxies” issue

- ◆ The proxy is only delegated if required
  - The condition is  $\text{lifetime} < 4 \text{ hours}$ .
- ◆ The delegation is performed by the `glite-transfer-submit` CLI. The first submit client that sees that the proxy needs to be redelegated is the one that does it - the proxy then stays on the server for ~8 hours or so
  - Default lifetime is 12 hours.
- We found a **race condition** in the **delegation** - if two clients (as is likely) detect at the same time that the proxy needs to be renewed, they both try to do it and this can result in the delegation requests being mixed up - so that that what finally ends up in the DB is the **certificate** from one request and the **key** from the other.
- ◆ We don't detect this and the proxy remains invalid for the next ~8 hours.
- ◆ The real fix requires a server side update (ongoing).

# BDII Scaling Problem

- ◆ BDII/SRM problem @ NIKHEF / SARA
- ◆ Discovery : only possible via monitoring of
  - Jobs success by exp'ts (not always optimum)
  - Site services by site
  - Coupled phenomenon
- ◆ BDII developer hears via 'vocal site person' about situation
- ◆ Active support
  - Checking deployment scenario
  - Asking for log files
  - Making recommendations

# Post mortem by developer

## Summary:

After fixing the initial problem with the missing index., the SRM failed at SARA and the lhcb jobs went into a lcg-gt loop which put a high query load onto the BDII.

Here are a number of recommendations following the incident.

Recommendation 1: Update the deployment and trouble shooting documentation explaining clearly the dangers of co-hosting the BDII with other services which could generate a high load.

Recommendation 2: Produce a new release of the BDII in which the new index is set.

Recommendation 3: Improve the efficiency of the lcg-utils commands, in particular has some kind of cache is need to avoid repeated queries to to BDII.

Recommendation 4: Use the logs gathered from the incident to test the performance of the BDII in such situation and address and performance bottlenecks found.

Recommendation 5: Monitor the load on the BDII at NIKHEF. If the load is consistently high, consider adding and additional machine for load balancing.

Recommendation 6: Advise the VOs the dangers of fail over methods that can do DOS loops. All fail over should contain some exponential backup.

Recommendation 7: Implement, as planned service discovery APIs. These interfaces to the information systems should contain limiter that prevent single threads (better processes) from issuing rapid fire queries. This limiter is needed to prevent accidental DOS attacks that make the whole resource unusable

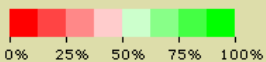
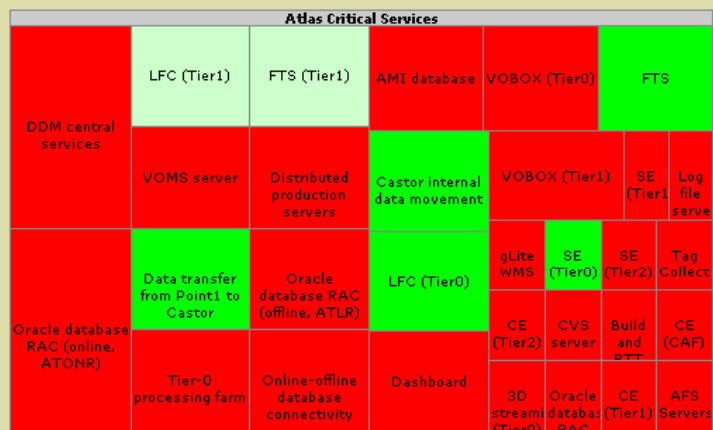


# CCRC Goals

- ◆ For February, we have the following three sets of metrics:
  1. The **scaling factors** published by the experiments for the various **functional blocks** that will be tested. These are monitored continuously by the experiments and reported on at least weekly;
  2. The lists of **Critical Services**, also defined by the experiments. These are complementary to the above and provide additional detail as well as service targets. It is a goal that all such services are handled in a standard fashion - i.e. as for other IT-supported services - with appropriate monitoring, procedures, alarms and so forth. Whilst there is no commitment to the problem-resolution targets - as short as 30 minutes in some cases - the follow-up on these services will be through the daily and weekly operations meetings;
  3. The services that a site must offer and the corresponding availability targets based on the **WLCG MoU**. These will also be tracked by the operations meetings.
- Phase 2 of CCRC in May

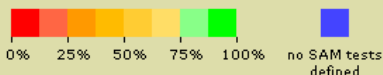
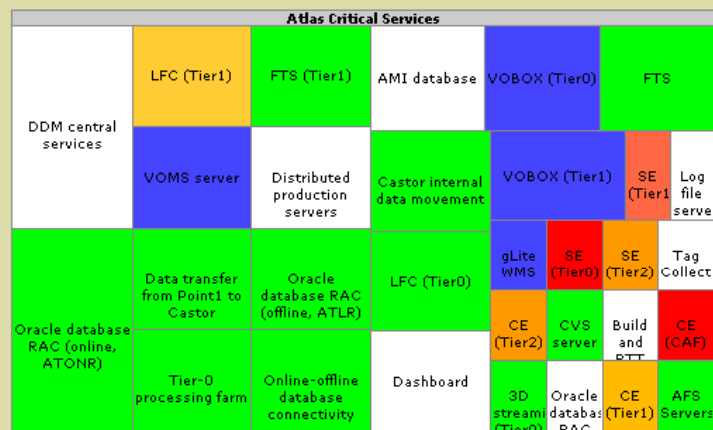
## WLCG CCRC'08 Critical Services "GridMap"

### Ticklist Status (updated manually)



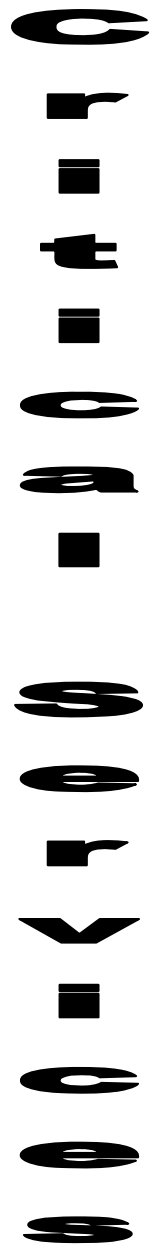
Alice **Atlas** CMS LHCb

### Test Status (live data)



CERN-PRDD TRIUMF-LCG2 IN2P3-CC FZK-LCG2 INFN-T1 SARA-MATRIX NDGF-T1 pic Taiwan-LCG2 RAL-LCG2  
BNL-LCG2 USCMS-FNAL-WC1

[bookmark settings](#)



# Service Readiness

| #  | Question   | Comments  |
|----|--|---|
| 1  | High-level description of service available?       | with architecture diagram                       |
| 2  | Middleware dependencies and versions defined?      | OS deps, M/W deps, platforms supported          |
| 3  | Code released and packaged correctly?              | Repository + Tagging process, rpms/tarballs     |
| 4  | Certification process exists?                      |   |
| 5  | Automatic Configuration code exists?               | e.g. Yaim, NCM, ...                             |
| 6  | Admin Guides available?                            | Installation, monitoring, problem determination |
| 7  | Disk, CPU, Database, Network requirements defined? |   |
| 8  | Monitoring criteria described?                     |   |
| 9  | Problem determination procedure documented         |   |
| 10 | Support chain defined (2nd/3rd level)?             |   |
| 11 | Backup/restore procedure defined?                  |   |
| 12 | Suitable hardware used                             |   |
| 13 | Monitoring implemented                             |   |
| 14 | Test environment exists                            |   |
| 15 | Problem determination procedure implemented        |   |
| 16 | Automatic configuration implemented                |   |
| 17 | Backup procedures implemented and tested           |   |

Key:

- Software Readiness
- Service Readiness
- Site Readiness

- Measure of how 'production-ready' a service :
  - In terms of software, service and deployment
- Manually edited (under SVN control) by responsables
  - EIS team, service managers, deployment team



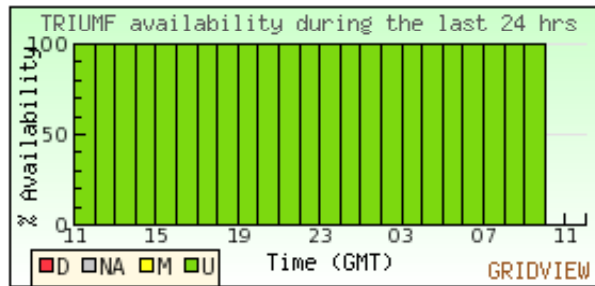
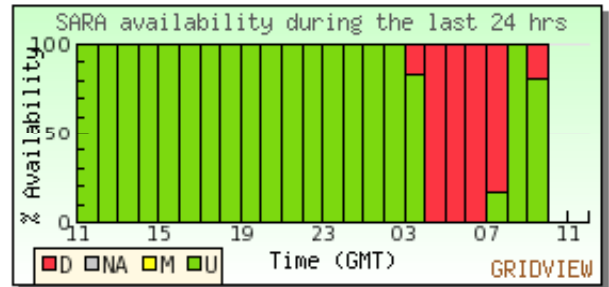
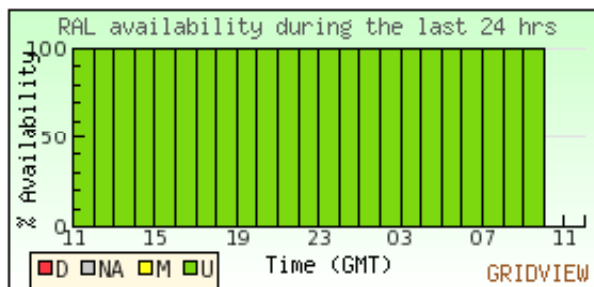
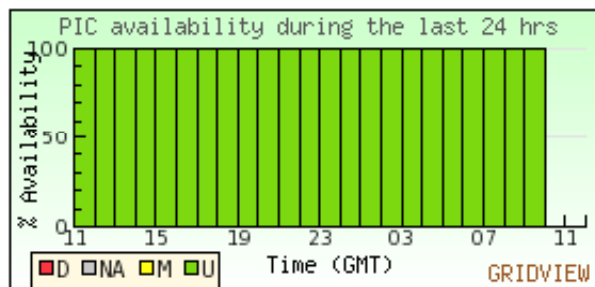
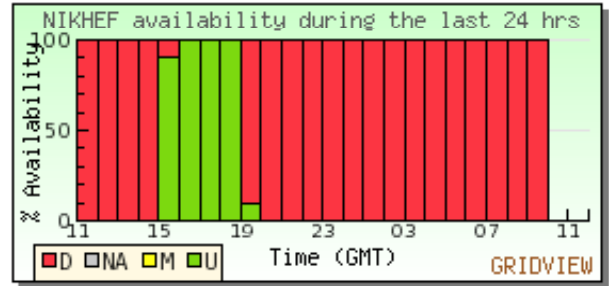
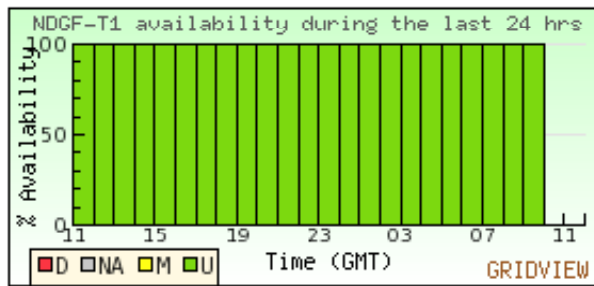
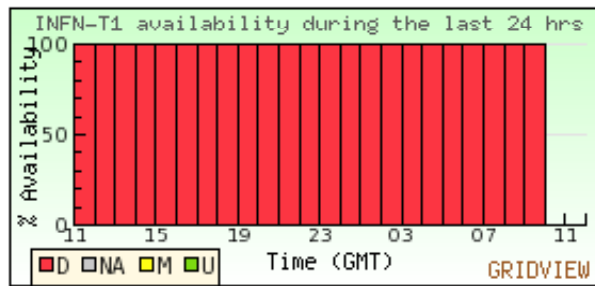
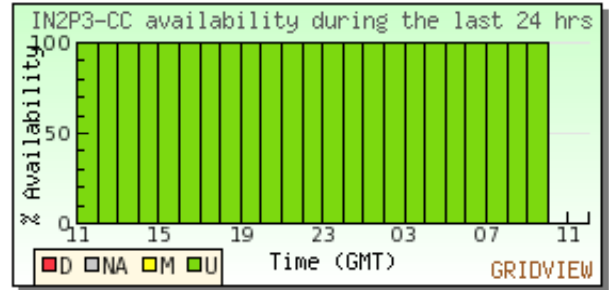
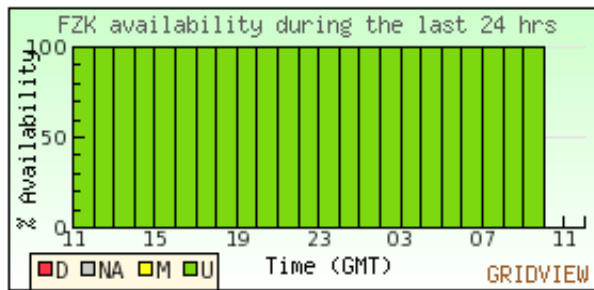
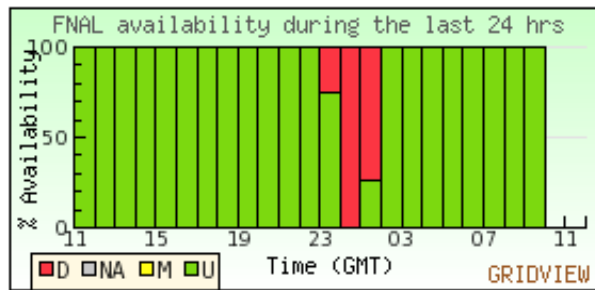
# CCRC Goals

- ◆ For February, we have the following three sets of metrics:
  1. The **scaling factors** published by the experiments for the various **functional blocks** that will be tested. These are monitored continuously by the experiments and reported on at least weekly;
  2. The lists of **Critical Services**, also defined by the experiments. These are complementary to the above and provide additional detail as well as service targets. It is a goal that all such services are handled in a standard fashion - i.e. as for other IT-supported services - with appropriate monitoring, procedures, alarms and so forth. Whilst there is no commitment to the problem-resolution targets - as short as 30 minutes in some cases - the follow-up on these services will be through the daily and weekly operations meetings;
  3. The services that a site must offer and the corresponding availability targets based on the **WLCG MoU**. These will also be tracked by the operations meetings.
- Phase 2 of CCRC in May

# Site performance

| <i>Service</i>   | <i>Maximum delay in responding to operational problems</i> |   |   | <i>Average availability measured on an annual basis</i> |                    |
|--|--|---|---|---|--------------------|
|  | Service interruption                                       | Degradation of the capacity of the service by more than 50% | Degradation of the capacity of the service by more than 20% | During accelerator operation                            | At all other times |
| Acceptance of data from the Tier-0 Centre during accelerator operation   | 12 hours   | 12 hours  | 24 hours  | 99%   | n/a                |
| Networking service to the Tier-0 Centre during accelerator operation   | 12 hours   | 24 hours  | 48 hours  | 98%   | n/a                |
| Data-intensive analysis services, including networking to Tier-0, Tier-1 Centres outwith accelerator operation | 24 hours   | 48 hours  | 48 hours  | n/a   | 98%                |
| All other services - prime service hours <sup>6</sup>  | 2 hour   | 2 hour  | 4 hours   | 98%   | 98%                |
| All other services - outwith prime service hours   | 24 hours   | 48 hours  | 48 hours  | 97%   | 97%                |

99% availability means < 100 minutes per week!



# GGUS is not fast enough ...

**Nagios**

**General**

- Home
- Documentation

**Monitoring**

- Tactical Overview
- Service Detail
- Host Detail
- Hostgroup Overview
- Hostgroup Summary
- Hostgroup Grid
- Servicegroup Overview
- Servicegroup Summary
- Servicegroup Grid
- Status Map
- 3-D Status Map
- Service Problems
- Host Problems
- Network Outages

Show Host:

- Comments
- Downtime
- Process Info
- Performance Info
- Scheduling Queue

**Reporting**

- Trends
- Availability
- Alert Histogram

**Current Network Status**  
 Last Updated: Thu Dec 6 15:37:42 CET 2007  
 Updated every 90 seconds  
 Nagios® - [www.nagios.org](http://www.nagios.org)  
 Logged in as /O=dutchgrid/O=users/O=nikhef/CN=Jeffrey Templon

**Host Status Totals**

| Up  | Down | Unreachable | Pending |
|-----|------|-------------|---------|
| 158 | 0    | 0           | 0       |

**Service Status Totals**

| Ok  | Warning | Unknown | Critical | Pending |
|-----|---------|---------|----------|---------|
| 856 | 6       | 0       | 10       | 3       |

**Display Filters:**  
 Host Status Types: All  
 Host Properties: Any  
 Service Status Types: All Problems  
 Service Properties: Any

**Service Status Details For All Hosts**

| Host                                | Service                                       | Status   | Last Check          | Duration      | Attempt | Status   |
|-------------------------------------|---|----------|---------------------|---------------|---------|--|
| <a href="#">bosheks.nikhef.nl</a>   | <a href="#">hr.srce.ResourceBroker-RunJob</a> | CRITICAL | 2007-12-06 15:29:43 | 0d 1h 15m 38s | 4/4     | Job subm failed:****<br>API_NATI<br>**** Error v<br>the NSCLie<br>api<br>Authentic<br>Failed to e<br>security co<br>Error:<br>UI_NO_N:<br>**** Unabl<br>any Netw |
| <a href="#">hooibroei.nikhef.nl</a> | <a href="#">check ncd</a>                     | WARNING  | 2007-12-06 15:33:34 | 1d 1h 34m 47s | 3/3     | WARNING<br>finished w<br>warning(s)  |
| <a href="#">hooikist.nikhef.nl</a>  | <a href="#">check ncd</a>                     | WARNING  | 2007-12-06 15:33:34 | 1d 1h 34m 47s | 3/3     | WARNING<br>finished w<br>warning(s)  |
| <a href="#">hooikuil.nikhef.nl</a>  | <a href="#">check ncd</a>                     | WARNING  | 2007-12-06 15:33:35 | 1d 1h 34m 46s | 3/3     | WARNING<br>finished w<br>warning(s)  |



# Middleware Coverage

- ◆ AAA : already reasonably well stressed
- ◆ FTS : broader range of usage, target SRMs, new SRM interface, higher rate (race condition ...)
- ◆ SRMs : stressed to max
- ◆ WMS : unknown to what extent. LHCb apparently using RB.
- ◆ CREAM : big miss. Should push extremely hard to get this ready for phase 2.
- ◆ Glexec / LCMAPS-server : another big miss
- ◆ All products could use improvement in logging / diagnostics / monitoring!!!!

# Handling Problems...

- ◆ Need to clarify current procedures for handling problems - some mismatch of expectations with reality
  - e.g. no GGUS TPMs on weekends / holidays / nights...
    - c.f. problem submitted with max. priority at 18:34 on Friday...
  - Use of on-call services & expert call out as appropriate
    - {alice-,atlas-}grid-alarm; {cms-,lhcb-}operator-alarm;
  - Contacts are needed on all sides - sites, services & experiments
    - e.g. who do we call in case of problems?
- ◆ Complete & open reporting in case of problems is essential!
  - Only this way can we learn and improve!
  - **It should not require Columbo to figure out what happened...**
- ◆ Trigger post-mortems when MoU targets not met
  - This should be a light-weight operation that clarifies what happened and identifies what needs to be improved for the future
  - Once again, the problem is at least partly about communication!

# Don't panic

- ◆ Many EGEE / JRA1 services are in considerably better shape than exp't middleware
- ◆ BUT this is no license to slow down or slack off :
  - exp't efforts are often much more focused, they can catch up quickly
  - EGEE services are more critical : problems here affect all VOs / entire site. You *\*must\** do better!
  - If exp'ts catch up and pass us, they will be merciless