



ATC-ABCO Days

Session 4 - MTTR & Spare Policy for the LHC injectors & experimental areas: AT & IT Groups

Databases, Networks, Informatics
22 January 2008

Tim Smith, Frédéric Hemmer

Abstract / Contents

The presentation will focus on describing IT services believed to be essential for accelerator operations, in particular for databases and networks, and how critical service(s) achieve high availability, as well as what level of coverage and/or standby service are available throughout the year.

A review of recent or significant incidents will be presented as well as measures taken to avoid or improve service recovery.

Finally a number of significant outstanding issues will be presented.



Computer Centre Operations

- 24x365 Operator on shift
 - Performs simple documented interventions
- 24x365 System Administration Coverage
 - Most of IT servers, incl. Linux DB servers
 - First line diagnosis & intervention
 - Answer within minutes; on site < 1 hour
 - Unyielding problems are forwarded to experts
- Experts on best effort coverage
 - Usually complex services
 - Most services do not have enough people to provide a piquet service

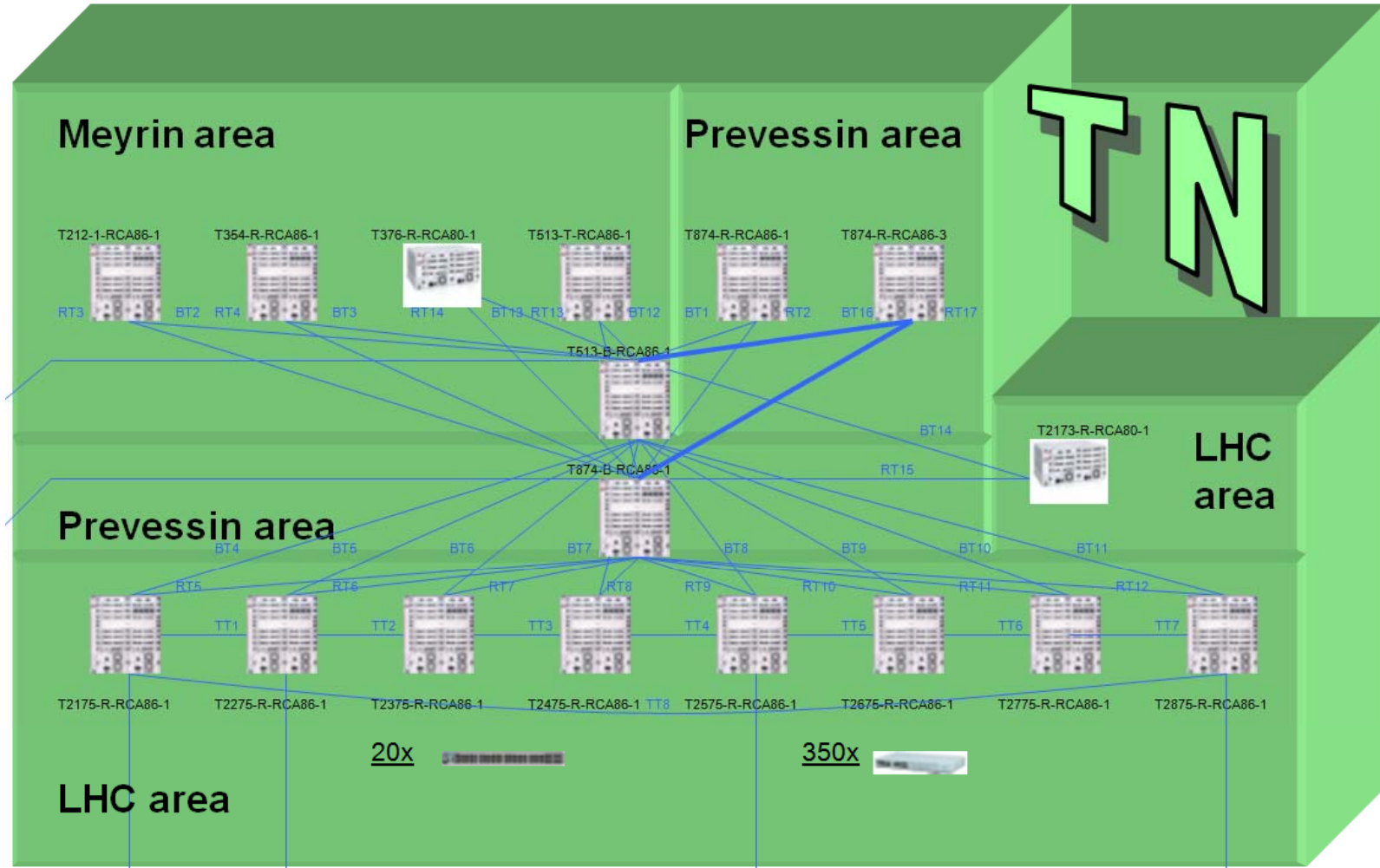


Communication services

- Cabling on the Technical & GPN networks
 - Guaranteed 20 years
 - Characteristics recorded
 - Guaranteed to work if untouched
- Equipment on Technical & GPN networks
 - Guaranteed for 5 years
 - 24x365 maintenance for replacement
 - Maximum 1 hour intervention time outside working hours
 - Stock available except for major disasters

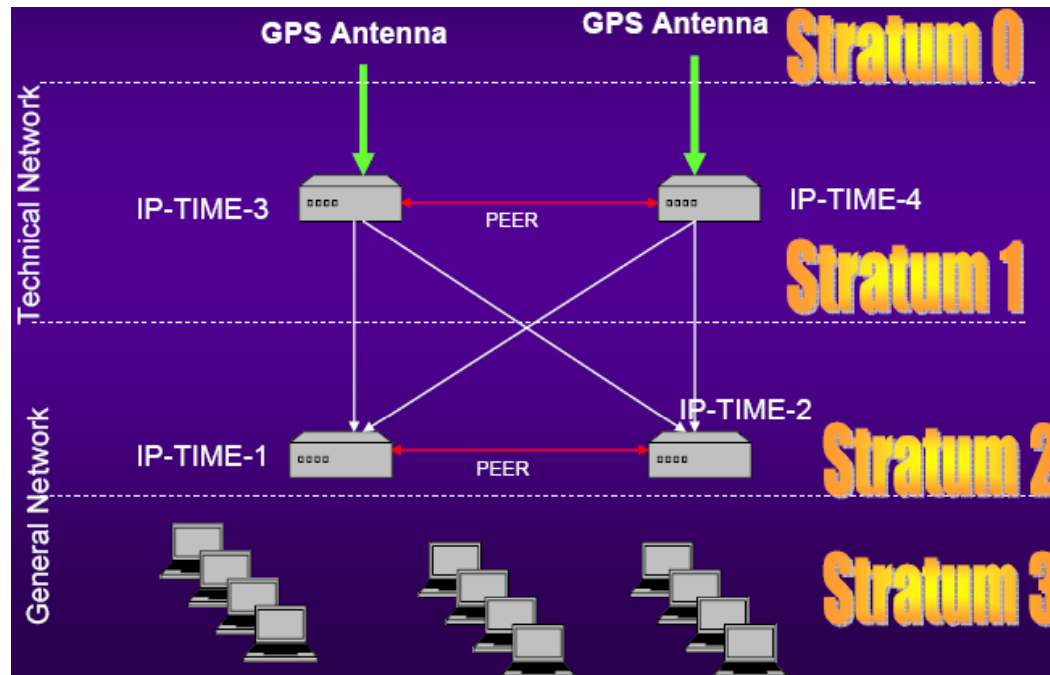


Campus Network Backbone



Communication services (II)

- Networking *Services*
 - DNS, DHCP, NTP, RADIUS
 - Redundant configurations
 - Best effort support outside worki



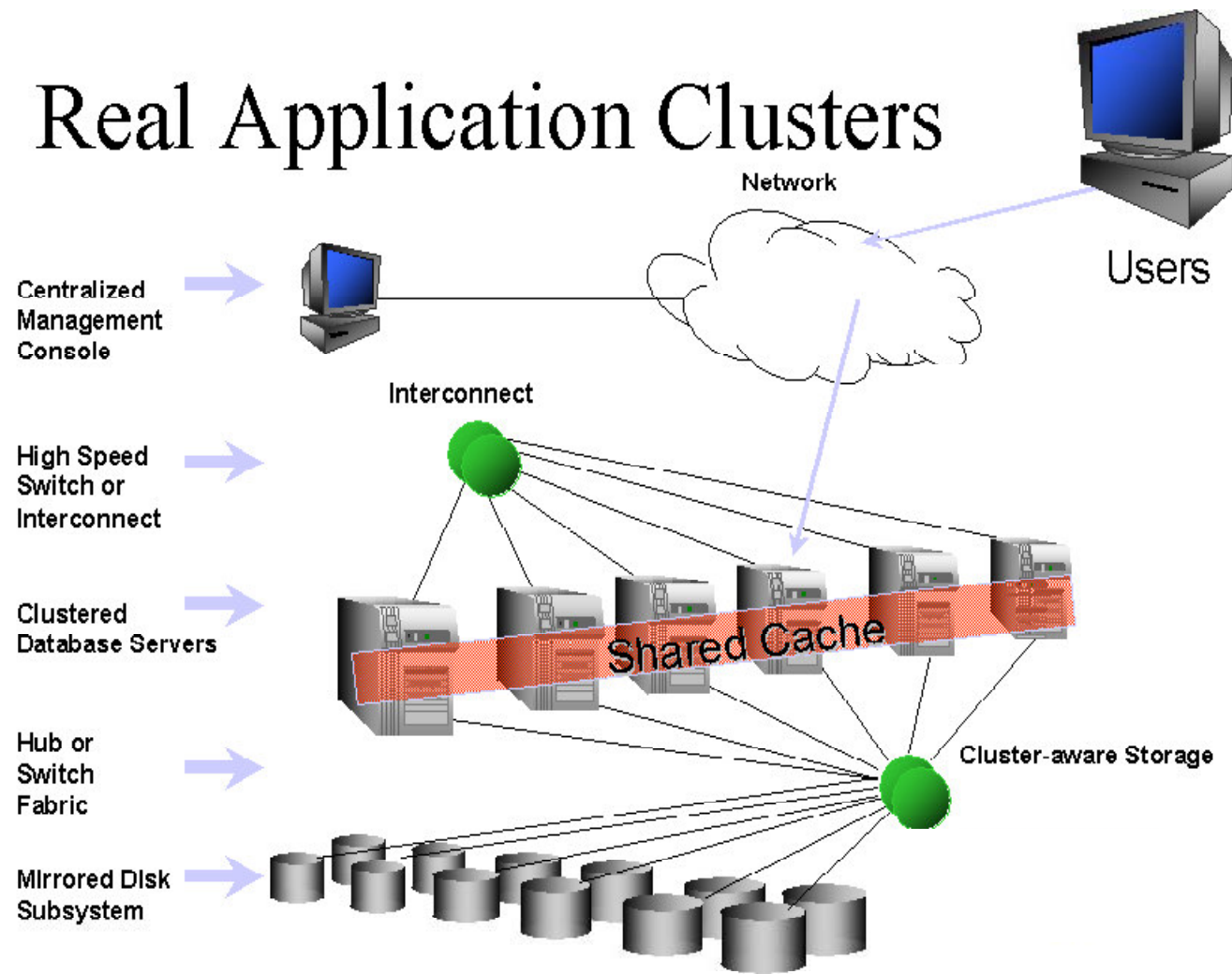
- Telephony Services
 - Redundant configurations
 - Maintenance contract 24x365 for fixed and mobile telephony (NextiraOne/Sunrise)
 - Max. 1hr intervention time outside working hours
 - The fixed telephone network has 3 hours of power autonomy
 - Local UPS (+ diesel backup...)
 - The GSM network is not covered by UPS



- Most recent DB services have higher redundancy
 - Linux with Oracle RAC
 - Storage with high availability features
 - Basic Server Administration using Computer Center Operation Services
 - Complex problems need experts to intervene
 - Best effort coverage
 - Most experts > Eb hence no compensation possible
- Some Databases are still running on legacy ageing Sun/Solaris
 - Planned to be upgraded to Linux/RAC in 1H08



Real Application Clusters



Significant incidents (I)



- *lhlogdb* currently serves accelerator settings, measurements and logging, all in one database
 - Recent major *lhlogdb* incidents
 - 15 Jan 2008: accidental data deletion, long restore of a copy
 - 27 Nov 2007: internal database bug hit, complete cluster stop, five minutes downtime, patch installed at next maintenance period
 - 8 Nov 2007: one system disk “semi-failure” leads to cluster instability and four hours downtime
 - 5 Jan 2006: power cut / Sun disk array cache issue and corruption, long restore/recover
- New platform for AB databases:
 - End 2007: new HW purchased; Linux RAC, implementing DataGuard/standby; being installed
 - “Settings” database will be split from “logging”
 - Several copies, smaller database (much faster restore/recover if required)

Significant incidents (II)



- An example of the exposure of old systems... major EDH database incident 2-5 Oct 2007
 - 3 days unavailability
 - Root cause was triple disk failure, data corruption
 - Very long restore/recover operation (no data loss)
 - Ageing platform without all the new storage features:
 - double parity, automated regular media check, checksum at the storage level
- Migration project has started
 - New platform is the same as the new AB databases server/storage
 - New platform will have additional copy of backup on disk and increased bandwidth to/from the tape backup system

Significant incidents (III)



- Complete power failure July 2006
 - Communication is the key factor
 - Rapid communication of diesel failure would have enabled pro-active shutdown of critical services
 - facilitating service restoration (and avoiding hardware failures?)
 - Interventions teams (elec, hvac) didn't contact IT operators
 - Over 30 experts came in and rapidly restored services
 - Despite lack of any formal piquet arrangements
 - Service redundancy can be invalidated by client configuration
 - One AB *service* restored without failed HW
 - But CCC clients were configured to use failed HW directly
 - Much confusion – *service* available, but clients cant use it

Significant incidents (IV)

- During the last general power cut
 - Two PBX servers went down
 - 8 hours, diesel did not start correctly
 - The complete GSM network went down
 - *Not on UPS*
- During the emergency power tests
 - The GSM network in the LHC stopped
 - *Not on UPS*



Issues (as perceived from IT)

- Need list of IT services critical for accelerator operations
 - The experiments have such a list including
 - Criticality; Responsible
 - Maximum allowed down time
 - Impact on the experiment
 - See <https://twiki.cern.ch/twiki/bin/view/CMS/SWIntCMSServices>
 - Implementation must take account of global resource envelope
- Interdependence Tests should be made by switching off IT services (or access to them);
 - Reliable services hide possible failure modes
- Databases need to be regularly updated with the quarterly Oracle security patch and relevant OS patches
- LHC “logging” database seems to be critical
 - Maybe could consider hardening applications by caching data?



Issues (as perceived from IT) - II



- Providing coverage better than “best effort” for IT services is problematic
 - Modern services are complex
 - Complicated end-to end problems require experts
 - Most services do not have the minimal number of experts required for standby services
 - People will not be willing to enroll to standby services if they are not compensated appropriately
 - Most of the experts are > Eb
 - IT services run the whole year
 - This problem has been highlighted for the last 7 years



- How to effectively communicate notice of service changes or interruptions
 - TS use *notes du coupure*
 - Printed and posted to entrance doors
 - IT use pop-up alerts targeted to impacted community
 - Coordinated through a *Service Status Board*
 - Which communities need to know about which changes / interruptions?
 - Returns to the criticality / dependency issue

CNIC Security Issues

- CNIC policy and implementation documents updated
 - Non-implementation leaves serious risks... *However*
 - Currently impossible because of commissioning
 - Then will be impossible because of operations
 - Then will be impossible because of machine development ...
This is living dangerously!
- Technical Network security is compromised by the significant number of “trusted” hosts
 - Especially important are desktop Development PCs
- TN Intrusion Detection System should be implemented
- Security of Controls PCs should be assessed/improved
- Reorganize connectivity of controls devices
- Regular security scans must be scheduled on the TN
- Review and reduce number of service accounts...



Summary

- IT Services are critical to accelerator operations
 - As illustrated in recent incidents
- Interdependencies are either unknown or undocumented
 - A list of critical services should be established
 - Tests should be performed to expose the dependencies
- 24x365 coverage applies to first line interventions only
 - Complex problems require experts who are only available on a best effort basis
- There are significant security risks with devices connected to the technical network
 - CNIC policy should be implemented
 - Regular scans and updates are necessary

