

PSS

Physics Services Support

CERN IT  
Department

# Physics Database Services at CERN

[physics-database.support@cern.ch](mailto:physics-database.support@cern.ch)

Maria Girone, CERN IT-PSS  
WLCG Tier2 Tutorials, CERN, June 2006



- How to build a reliable and redundant database service?
  - Hardware choices
  - Procedures
- What role does the database service have in WLCG?

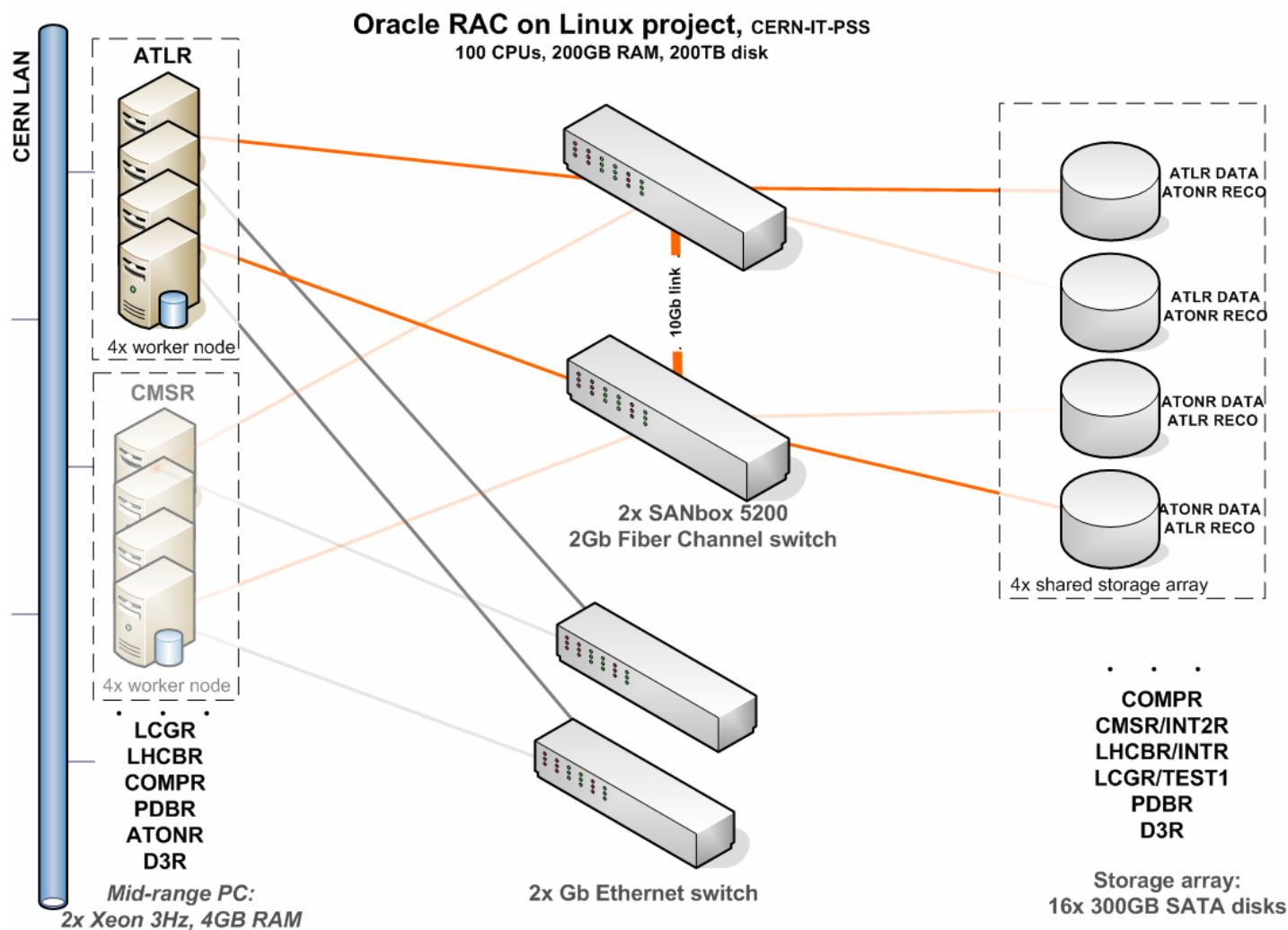
- Oracle services at CERN Tier 0 are used for
  - Conditions data
  - File Transfers
  - File Catalogs
  - Castor
  - Other experiment and Grid Applications
    - bookkeeping, physics production processing, on-line integration, detector construction and calibration, grid monitoring
- Database distribution outside Tier 0 are handled by the 3D project
  - ORACLE at Tier 1
  - Possibly mysql at Tier2
  - CERN tier0 is one participating site
  - More info at [lcg3d.cern.ch](http://lcg3d.cern.ch)

- Oracle services at CERN Tier 0 are used for
  - Conditions data
  - File Transfers
  - File Catalogs
  - Castor
  - Other experiment and Grid Applications
    - bookkeeping, physics production processing, on-line integration, detector construction and calibration, grid monitoring
- Database distribution outside Tier 0 are handled by the 3D project
  - ORACLE at Tier 1
  - Possibly mysql at Tier2
  - CERN tier0 is one participating site
  - More info at [lcg3d.cern.ch](http://lcg3d.cern.ch)

- Mandate: offer a **highly available** and **scalable** database service to the LHC experiments and grid deployment teams
  - Scalability - in both database processing power and storage
  - Flexibility - to cope with increasing demand
  - Reliability - automatic failover in case of problems
  - Manageability - significantly easier to administer than many individual disk servers
  - Isolation - 10g 'services' and/or physical separation
- Architecture choice
  - Database software -> Real Application Cluster 10g
  - Operating system -> Linux (RedHat ES)

- Summer 2005
  - Solaris based shared Physics DB cluster (2-nodes for HA)
    - Low CPU power, hard to extend, shared by all experiments
  - 40 (many) linux disk servers as DB servers
    - High maintenance load, no resource sharing, no redundancy
- Autumn 2005: consolidation on extensible database clusters (RAC)
  - No sharing across experiments
  - Higher quality building blocks
    - Midrange PCs (RedHat ES)
  - FibreChannel attached disk arrays
- Hardware resources more than doubled, same DBA team

- The Physics Database Services are deployed on 4-node and 2-node RAC/Linux, in failover mode



- Linear ramp-up budgeted for hardware resources in 2006-2008
- Planning next major service extension for Q3 this year (current resources will be doubled)

Current state (summer 2006)						
ATLAS	CMS	LHCb	Grid	3D	Non-LHC	PDB
4-node	4-node	4-node	4-node	2-node	4-node Compass	2-node
2-node valid/test	2-node valid/test	2-node valid/test	2-node pilot			
2-node online test						



## Service Size

- 50 mid-range servers and ~50 disk arrays (~600 disks)
- In other words: 100 CPUs, 200GB of RAM, 200 TB of raw disk space
- Half of the servers are in production, monitored 24x7
- ORACLE 10gR2 as main platform

## Service Procedures

- On-call team for 24x7 coverage
  - 4 DBAs and 5 developers (2 people on call)
- Backups on tape and on disk
- Recovery procedures validated
  - Default backup retention policy and frequency to be agreed with experiments/projects
- Monitoring: Oracle Enterprise Manager for DBAs
  - Application monitoring for users being integrated in Lemon

### Development Service

- Code development, no large data volumes, no backups
- one shared cluster
- 8x5 monitoring and availability

### Validation Service

- Larger tests and optimization
- 2-node RAC clusters
- 8x5 monitoring and availability
- DBA consultancy

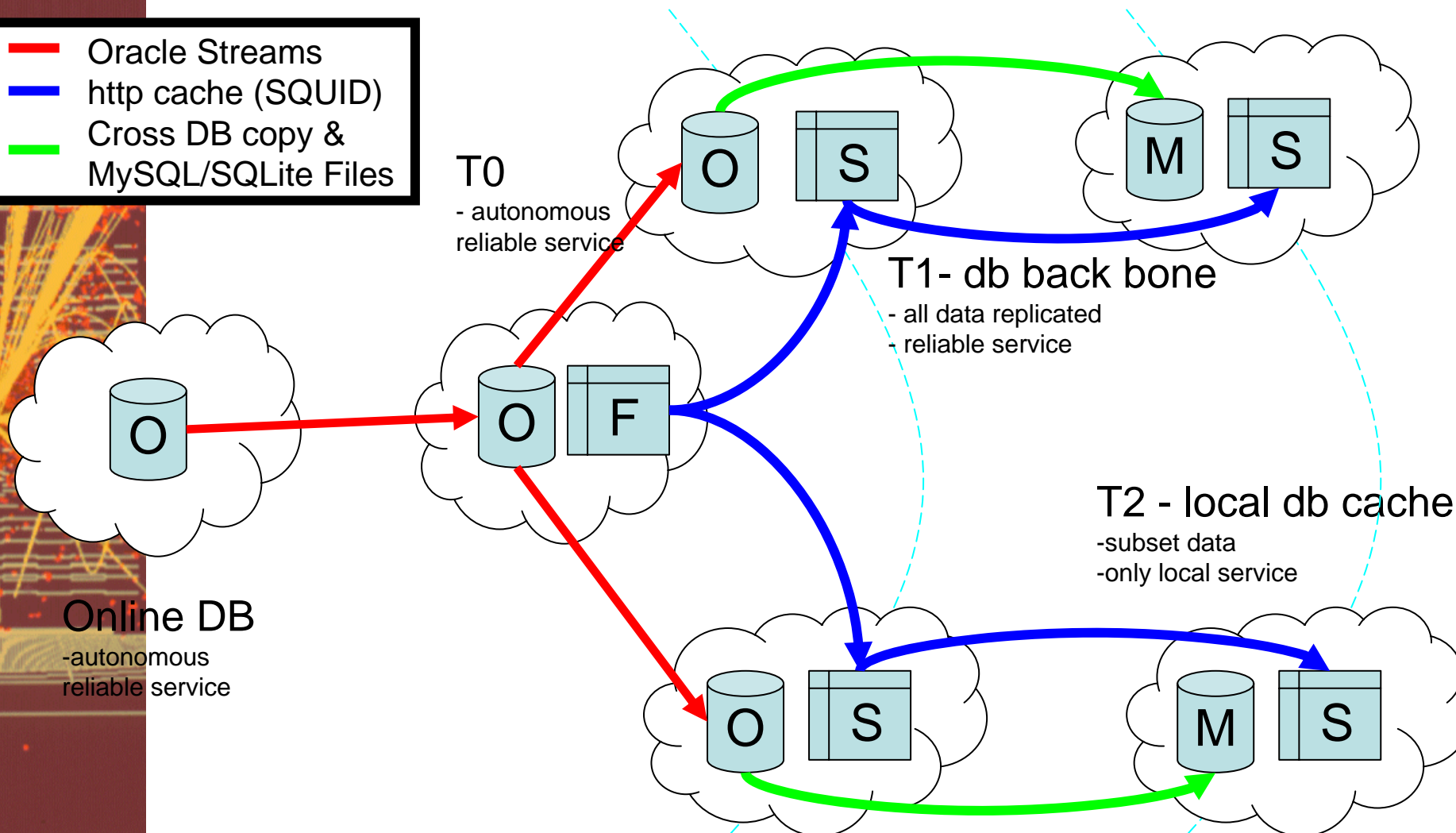
### Production Service

- 24x7 monitoring and availability, on call intervention procedures
- 4-node RAC cluster
- Backups every 30 minutes
- Limited number and scheduled planned interventions

- Resource usage report to experiment and project database coordinator
  - Allow experiment to prioritize resources and identify unexpected usage patterns
  - Which jobs/users got affected by what limit?
- Resource allocation and planning done together with the experiments, using these reports

- Oracle services at CERN Tier 0 are used for
  - Conditions data
  - File Transfers
  - File Catalogs
  - Castor
  - Other experiment and Grid Applications
    - bookkeeping, physics production processing, on-line integration, detector construction and calibration, grid monitoring
- Database distribution outside Tier 0 are handled by the 3D project
  - ORACLE at Tier 1
  - Possibly mysql at Tier2
  - CERN tier0 is one participating site
  - More info at [lcg3d.cern.ch](http://lcg3d.cern.ch)

- Oracle Streams
- http cache (SQUID)
- Cross DB copy & MySQL/SQLite Files



Online DB  
-autonomous  
reliable service

T0  
- autonomous  
reliable service

T1- db back bone  
- all data replicated  
- reliable service

T2 - local db cache  
-subset data  
-only local service

R/O Access at Tier 1/2  
(at least initially)

Dirk Duellmann, CERN IT

- Physics Database services fully based on RAC
  - Benefits of consolidation and additional flexibility obtained
- We have achieved a highly available and scalable service
  - We are ready for the challenges of the LHC start-up
- Q3 Database extension planned
  - The database resources will be doubled again