



# Considerations for database servers

Castor review – June 2006

Eric Grancher, Nilo Segura Chinchilla IT-DES



# Outline



## ❖ Current implementation

- Machine setup
- Database setup
- Throttling

## ❖ Return on experience

## ❖ Documentation, operations and sharing with Tier1 sites

## ❖ Upcoming implementation

## ❖ Our conclusions



# Questions to be addressed



- ❖ Only for the database part
- ❖ 3. Is the service infrastructure, and particularly the database server hardware, appropriate?
- ❖ 4. Are the necessary operational procedures in place (with adequate trained personnel) to ensure the MoU service reliability commitments can be met?
- ❖ 5. Does the overall service have the required performance and scalability the Tier0? For the CAF? What aspects of the software or system limit performance or scalability? Are these limits a concern?
- ❖ 6. Is the current policy for the support of Castor at other sites effective and sustainable?



# Return on the database architecture



- ❖ Castor(1) name server database
- ❖ Castor(2) stager database
- ❖ Castor(2) distributed logging facility database
- ❖ SRM database
  
- ❖ Different workload, can be grouped.



# Current implementation



- ❖ 1 Castor name server database
- ❖ 6 Castor2 environments (each 2 database servers stager + DLF): ALICE, ATLAS, CMS, LHCb, ITDC, SC4
- ❖ Castor2 test (RAC system and DLF) and dev
- ❖ Castor2 SRM database
- ❖ -> 18 database systems
  
- ❖ Using “CERN disk servers”
- ❖ RedHat Enterprise Linux (3ES x86)
- ❖ Oracle Enterprise Edition database (10.2)



# Issues with “CERN disk servers”



## ❖ Hardware problems

- Lack of support from the vendor
- Hardware problems, corruptions (see next slide)

## ❖ Machines are not RedHat Enterprise certified

- Ultimately Oracle Support can reject calls (has happened: Oracle Support “ask your HW vendor to study the issue and review the certification”)

## ❖ Not designed for database workload

- Database systems require a lot of IO operations per second (typically 8kB) -the case for the Castor2 stager- and very few “large IO”

## ❖ Lack of performance / memory

- New preferred deployment Oracle platform is Linux x86-64
- Larger memory
- Several CPUs and/or cores



# Corruptions



- ❖ 11 corruptions worked on between November 2005 and April 2006
  - On control files
  - On redo log files
  - On data files

Fri Nov 11 01:16:20 2005

Errors in file /ORA/dbs00/oracle/admin/CASTORSG/bdump/castorsg\_arc0\_7289.trc:

ORA-00354: corrupt redo log block header

ORA-00353: log corruption near block 212974 change 1612488170 time 11/11/2005 01:00:59

- ❖ Corruptions typically imply service downtime (lengthy recovery) and lot of work, may generate data loss.
- ❖ For Castor2, no data lost (but data loss on similar types of hardware for other services).
- ❖ Many of these corruptions have been fixed by the DBA during the night, delicate to automate fix (depends on type / place of corruption, not doing the right thing might harm more than help!).
- ❖ (since 1996), never had any of these type of corruptions on other platforms (Oracle writes a block which can be read later but not consistent with the format).
- ❖ *Many thanks* to IT-FIO Castor deployment and Linux Support for their help / investigations / new machines installations on the hardware platform.



# Throttling



- ❖ The “database system” has to sustain the load and manage priorities (even in worst cases, where it should throttle the load while still being reactive)
- ❖ Throttling
  - CPU in case of resource starvation
  - Active sessions (better usage of resource, limits latch inter-locking)
- ❖ Strict limitations
  - PGA Memory (generate ORA-600 to be trapped)
  - Sessions (only for new session, no issue in stable usage, “garde-fou”)
  - SGA maxsize (should not be reached, but...)
  - Hints for the system
    - PGA target
    - SGA target
  - Space of course
    - Each user / major component has its own tablespace with limitations





# Castor NameServer (DB) scalability



- ❖ Using a test (thanks to Olof/Ben), multithreaded / loop
  - Cns\_creat
  - Cns\_stat
  - Cns\_setfsize
  - Cns\_filestat
- ❖ Without problem, up to 920 Castor nameserver operations per second (with 15 threads).
- ❖ Oracle database is fully instrumented. We use the “wait event interface” as the main source for performance tuning (and early indication of scalability issues).



# Castor NameServer (DB) figures



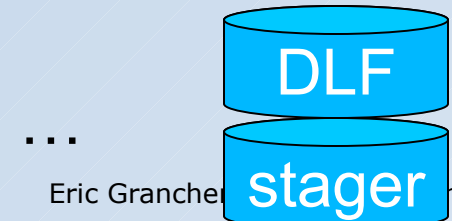
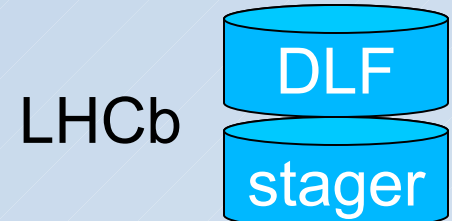
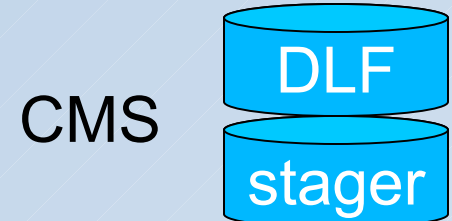
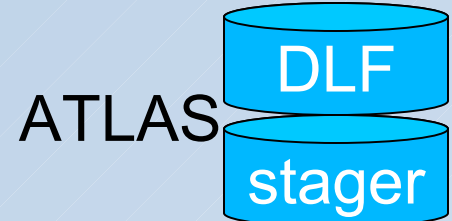
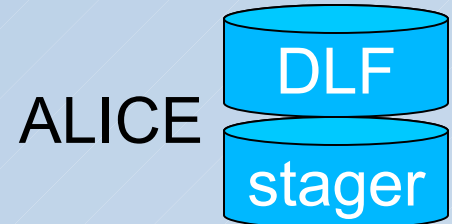
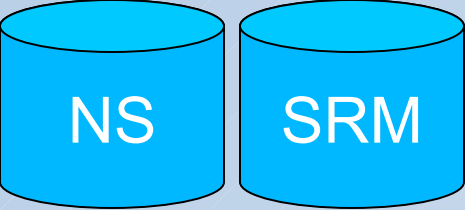
<b>Event</b>	<b>Waits</b>	<b>Time (s)</b>	<b>Avg Wait(ms)</b>	<b>% Total Call Time</b>	<b>Wait Class</b>
enq: TX - row lock contention	135,859	8,922	66	38.9	Application
log file parallel write	67,282	4,731	70	20.6	System I/O
log file sync	75,333	4,108	55	17.9	Commit
wait for scn ack	43,899	1,873	43	8.2	Other
control file sequential read	6,477	601	93	2.6	System I/O



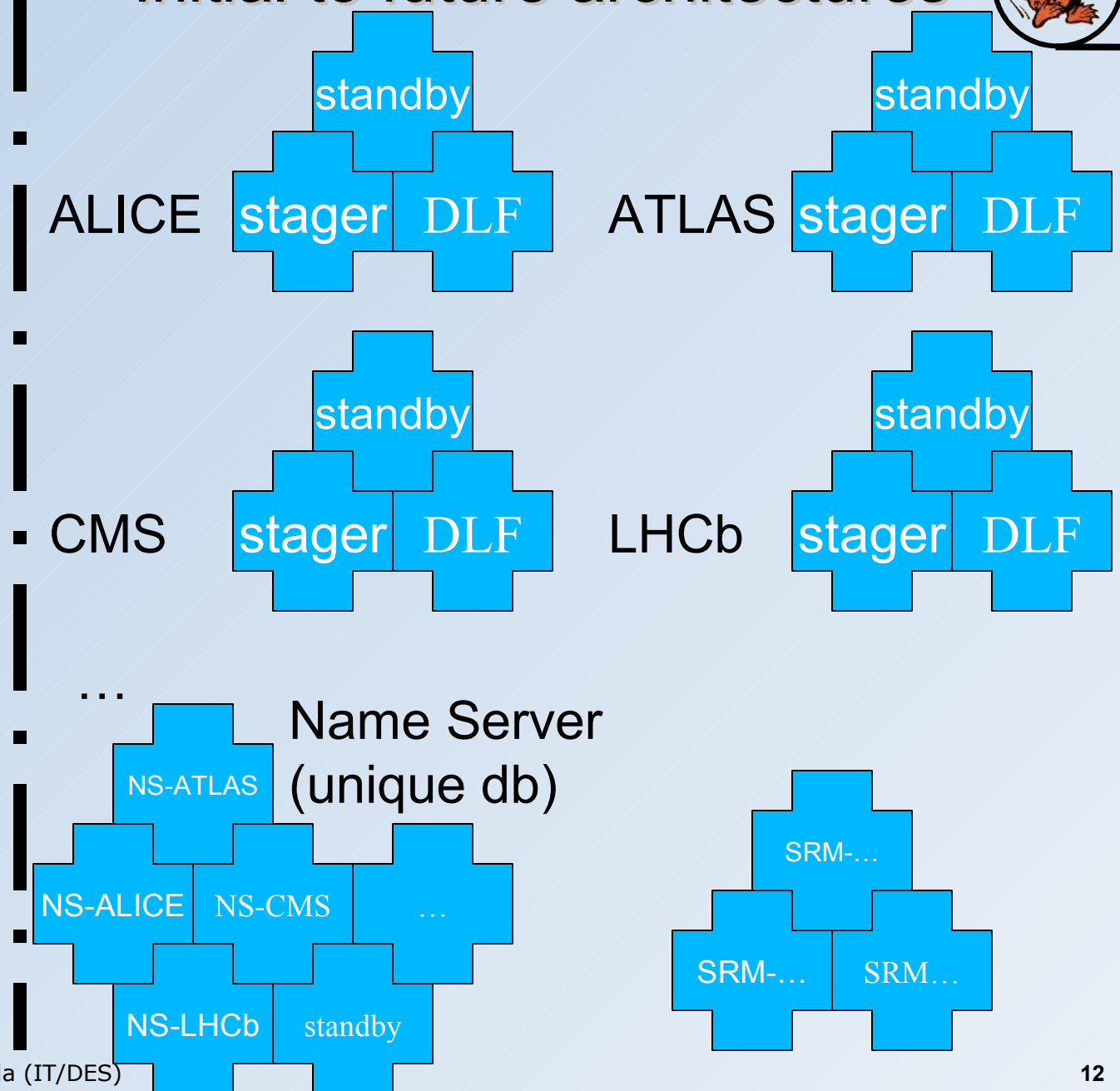
# New platform



- ❖ Use of the cluster (RAC) and standby (DataGuard) technologies understanding its limits
- ❖ Certified hardware (tested / validated / supported / guaranteed that the HW supplier and Oracle will work together on eventual issues)
- ❖ Flexible (depending on needs)
  - NS, stager, DLF, read-status
  - Eventually in a dynamic way
- ❖ Flexible (can have components made more powerful)
- ❖ Can scale with load



# Initial to future architectures





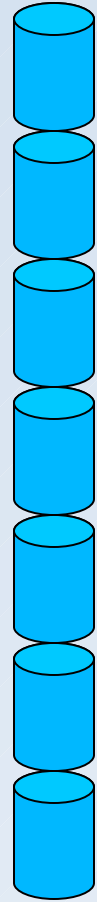
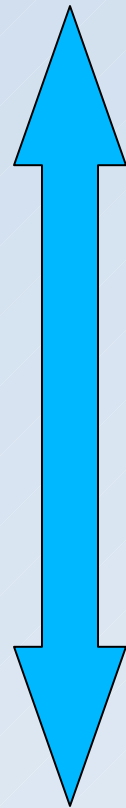
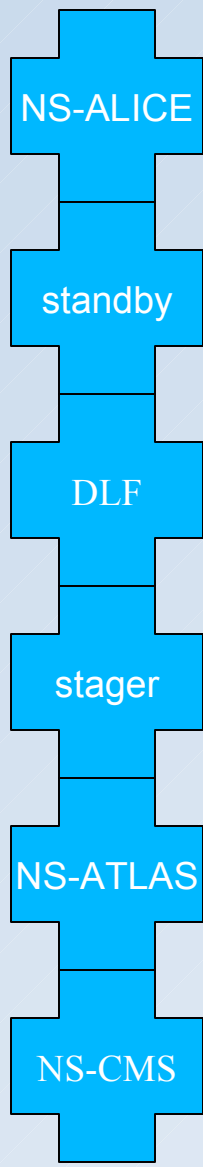
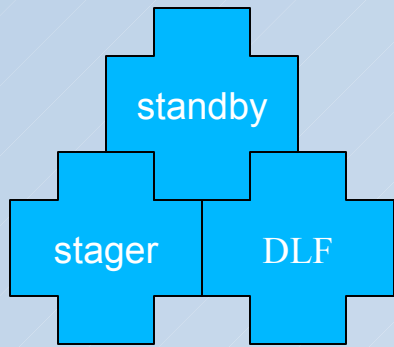
# Variations / combinations



- ❖ We may see that the usage patterns actually put a high-load on “name server”, stager, “status read”...
- ❖ Oracle clustering (RAC) make it possible to re-allocate resource where needed
  - 1 node for DLF + 2 nodes for stager versus
  - 1 node for stager + 2 nodes for DLF...
- ❖ Having standby (DataGuard) enables us to perform almost upgrades with very minimal downtime (<1 minute), cluster (RAC) does not.
- ❖ All database instances in a cluster share the same database, if there is a corruption, all nodes in the cluster may stop. With a standby, databases do not share “files”.
- ❖ Cluster and standby combination is the Oracle recommended maximum availability architecture.

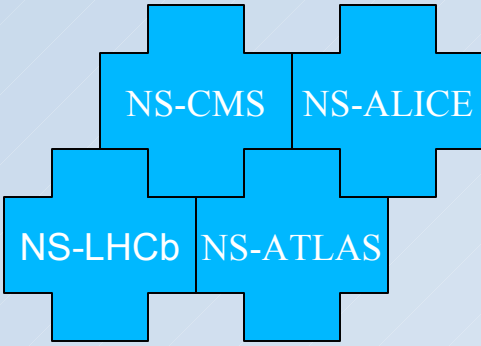


# Same hardware anyway



=

Name Server  
(unique db)



servers

storage network

storage arrays



# Collaboration with other institutes



- ❖ Oracle licensing.
  - Castor2 has been tested and works with Oracle Express Edition (free / limited version).
- ❖ Provide the creation scripts (with space and user allocation throttling/limits...) and can provide help for the initial deployment.
  - Full production support has to be done “locally”.
- ❖ Provide “Oracle statistics” (to have the same execution plans than the ones validated at CERN).
- ❖ Provide recipes / CERN Castor database guideline for deployment.
- ❖ Sharing of information (email for now, will to have phone calls, meetings).



# Operations



- ❖ Database administrators are part of castor-deployment mailing list.
  - Very good interaction with the Castor dev/deployment teams.
- ❖ Database administrators can be called 24x7 (we also cover Christmas) with rotation on the people. Operators have the procedure (and had to exercise it ☹, see corruptions).
- ❖ All operations are logged in Twiki (visible to Castor deployment team, experiments).
- ❖ Reference documents for installation
  - Space allocation
  - “Oracle user” privileges
  - Throttling and limits
  - Table statistics





# Security



## ❖ Model follows security best practices

- no extra processes ran on the database servers,
- No OS user login (administrators via ssh key).
- OS and Oracle patches applied on a regular basis).

## ❖ Use Oracle advanced security mechanisms can be used (encryption, more secure authentication)

- Tested and used on other services, non-significant overhead (especially for non-bandwidth intensive applications).
- CERN has Oracle license for it.
- To be validated, but a priori no issue to be faced.



# Our conclusions



- ❖ A lot of experience gathered with the “initial” deployment (for the throttling / limits, for the execution plans / together with the other sites).
- ❖ We have “encouraged” an early move for all databases to database version 10.2 (10gR2), “premier support” until July 2010 (extended 2013).
- ❖ Security of the platform is correct satisfactory, more can be done if required (transparent for the application).
- ❖ We believe, based on our experience, that adequate hardware / infrastructure is mandatory for a service to reach its quality and availability goals.



# References



- ❖ Oracle support

<http://www.oracle.com/support/library/data-sheet/oracle-li>

- ❖ Olof's stress test on the stager (stress on Oracle)

<http://castor.web.cern.ch/castor/presentations/2006/LSF->