

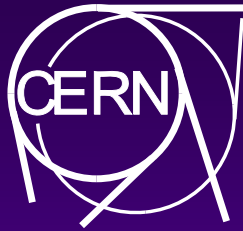
CERN local High Availability solutions and experiences

Thorsten Kleinwort

CERN IT/FIO

WLCG Tier 2 workshop CERN

16.06.2006



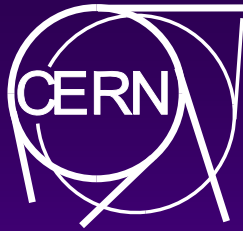
Introduction

- Different h/w used for GRID services
- Various techniques &
First experiences &
Recommendations



GRID Service issues

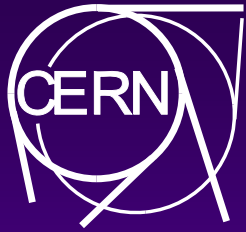
- Most GRID Services do not have built in failover functionality
- Stateful servers:
 - Information lost on crash
 - Domino effects on failures
 - Scalability issues: e.g. >1 CE



Different approaches

Server hardware setup:

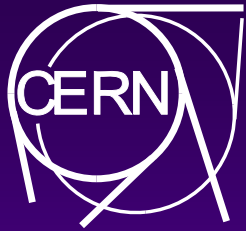
- Cheap 'off the shelf' h/w satisfies only for batch nodes (WN)
- Improve reliability (cost ^)
 - Dual Power supply (and test it works!)
 - UPS
 - Raid on system/data disks (hot swappable)
 - Remote console/BIOS access
 - Fiber Channel disk infrastructure
- CERN: 'Mid-Range-(Disk-)Server'



16 June 2006

Thorsten Kleinwort CERN-IT

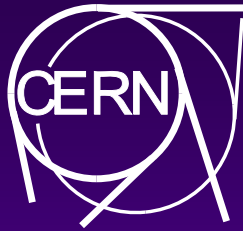
5



Various Techniques

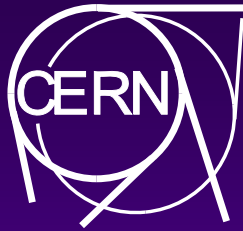
Where possible, increase number of 'servers' to improve 'service':

- WN, UI (trivial)
- BDII: possible, state information is very volatile and re-queried periodically
- CE, RB,... more difficult, but possible for scalability

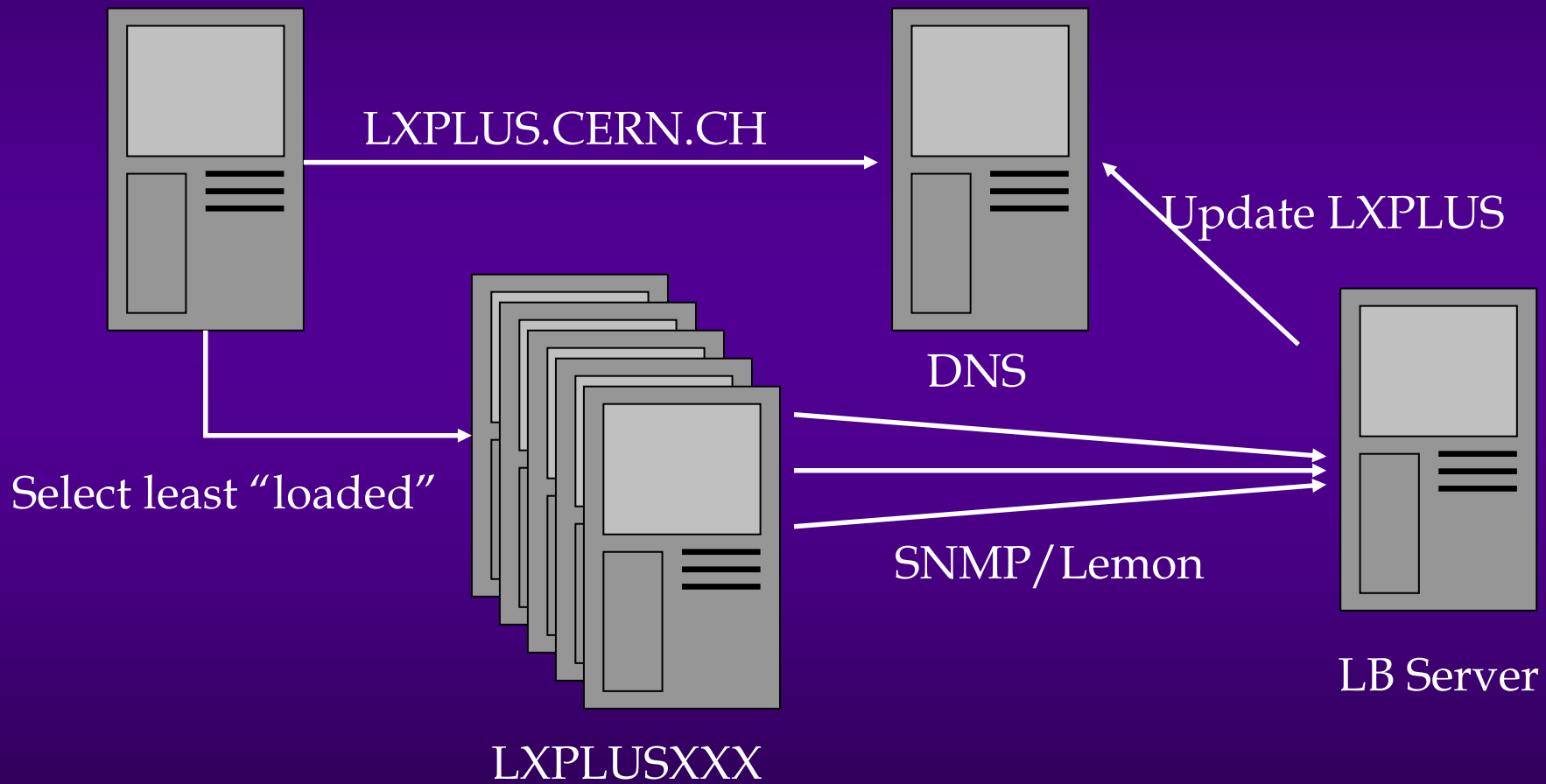


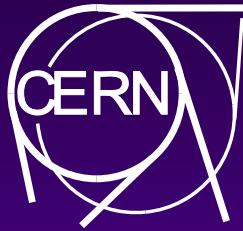
Load Balancing

- For multiple Servers, use DNS alias/load balancing to select the right one:
 - **Simple DNS aliasing:**
 - For selecting master or slave server
 - Useful for scheduled interventions
 - Dedicated (VO-specific) host names, pointing to the same machine (E.g. RB)
 - **Pure DNS Round Robin (no server check)**
configurable in DNS: to equally balance load
 - **Load Balancing Service, based on Server load/availability/checks:** Uses monitoring information to determine best machine(s)



DNS Load Balancing

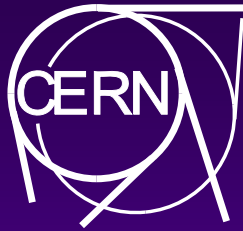




Load balancing & scaling

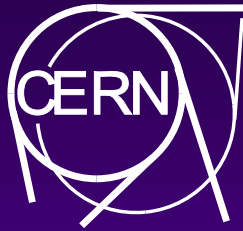
Publish several (distinguished) Servers for same site to distribute load:

- Can also be done for stateful services, because clients stick to the selected server
- Helps to reduce high load (CE, RB)



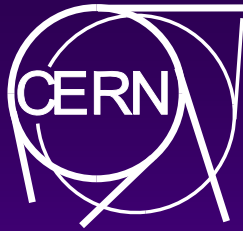
Use existing 'Services'

- All GRID applications that support ORACLE as a data backend (will) use the existing ORACLE Service at CERN for state information
- Allows for stateless and therefore load balanced application
- ORACLE Solution is based on RAC servers with FC attached disks

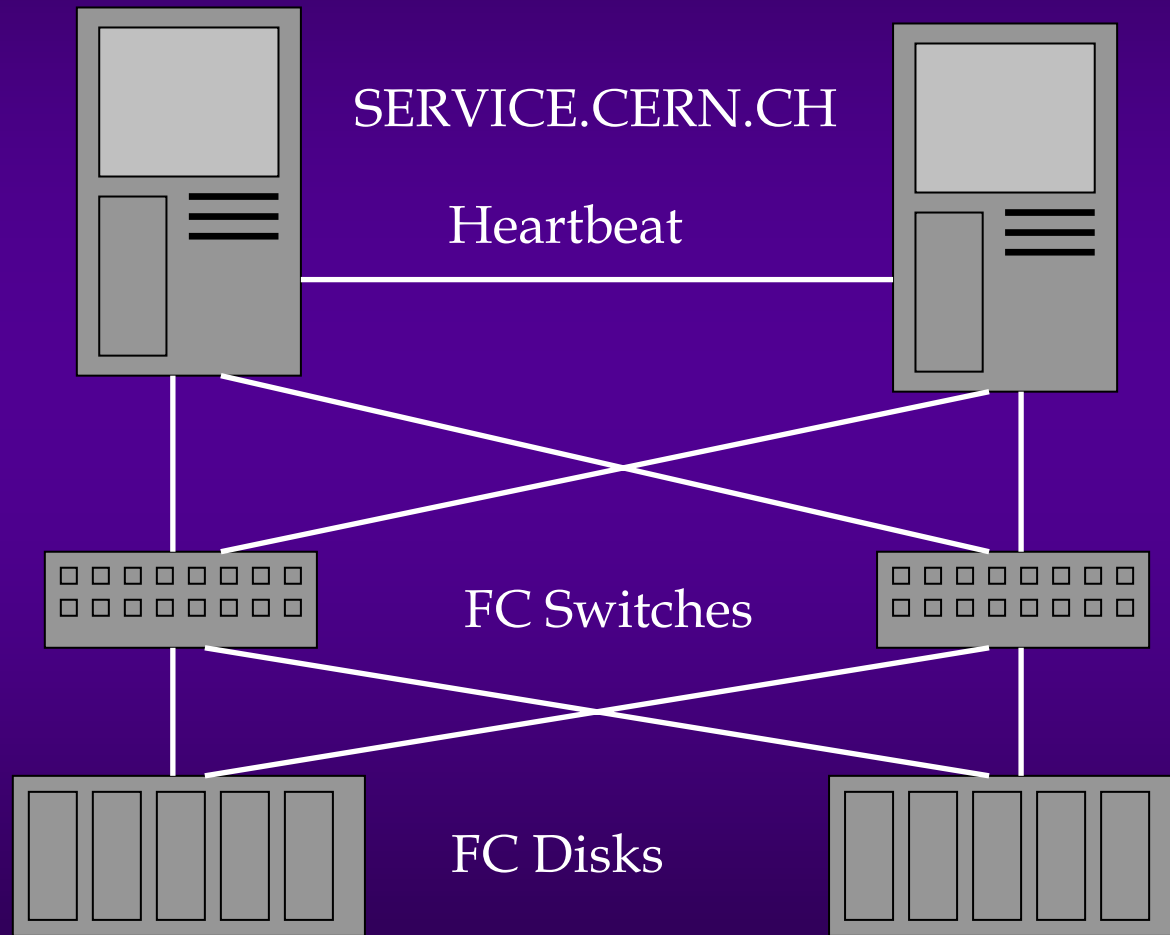


High Availability Linux

- Add-on to standard Linux
- Switches IP address between two servers:
 - If one Server crashes
 - On request
- Machines monitor each other
- To further increase high availability:
State information can be on (shared) FC



HA Linux + FC disk



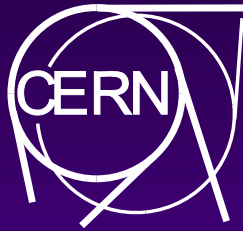
No single
Point of
failure



WMS: CE & RB

CE & RB are stateful Servers:

- load balancing only for scaling:
 - Done for RB and CE
- CE's and RB's on 'Mid-Range-Server' h/w
- If one server goes, the information it keeps is lost
- [Possibly FC storage backend behind load balanced, stateless servers]



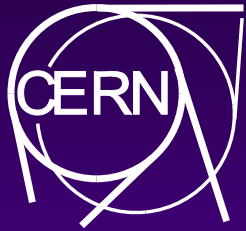
GRID Services: DMS & IS

DMS (Data Management System)

- SE: Storage Element
- FTS: File Transfer Service
- LFC: LCG File Catalogue

IS (Information System)

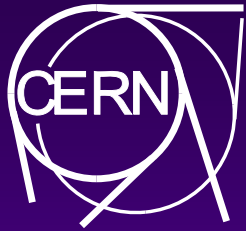
- BDII: Berkley Database Information Index



DMS: SE

SE: castorgrid:

- Load balanced front end cluster, with CASTOR storage backend
- 8 simple machines (batch type)
Here load balancing and failover works
(but not for simple GRIDFTP)



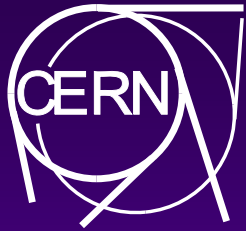
DMS:FTS & LFC

LFC & FTS Service:

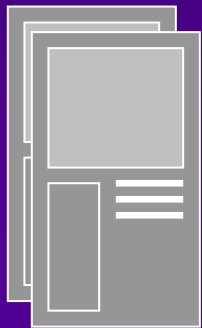
- Load Balanced Front End
- ORACLE database backend

+FTS Service:

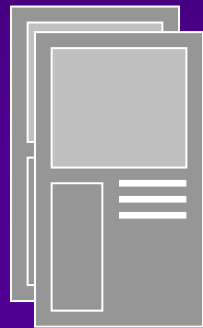
- VO agent daemons
One 'warm' spare for all:
Gets 'hot' when needed
- Channel agent daemons (Master & Slave)



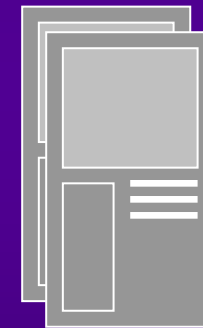
FTS: VO agent (failover)



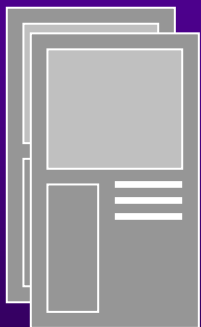
ALICE



ATLAS



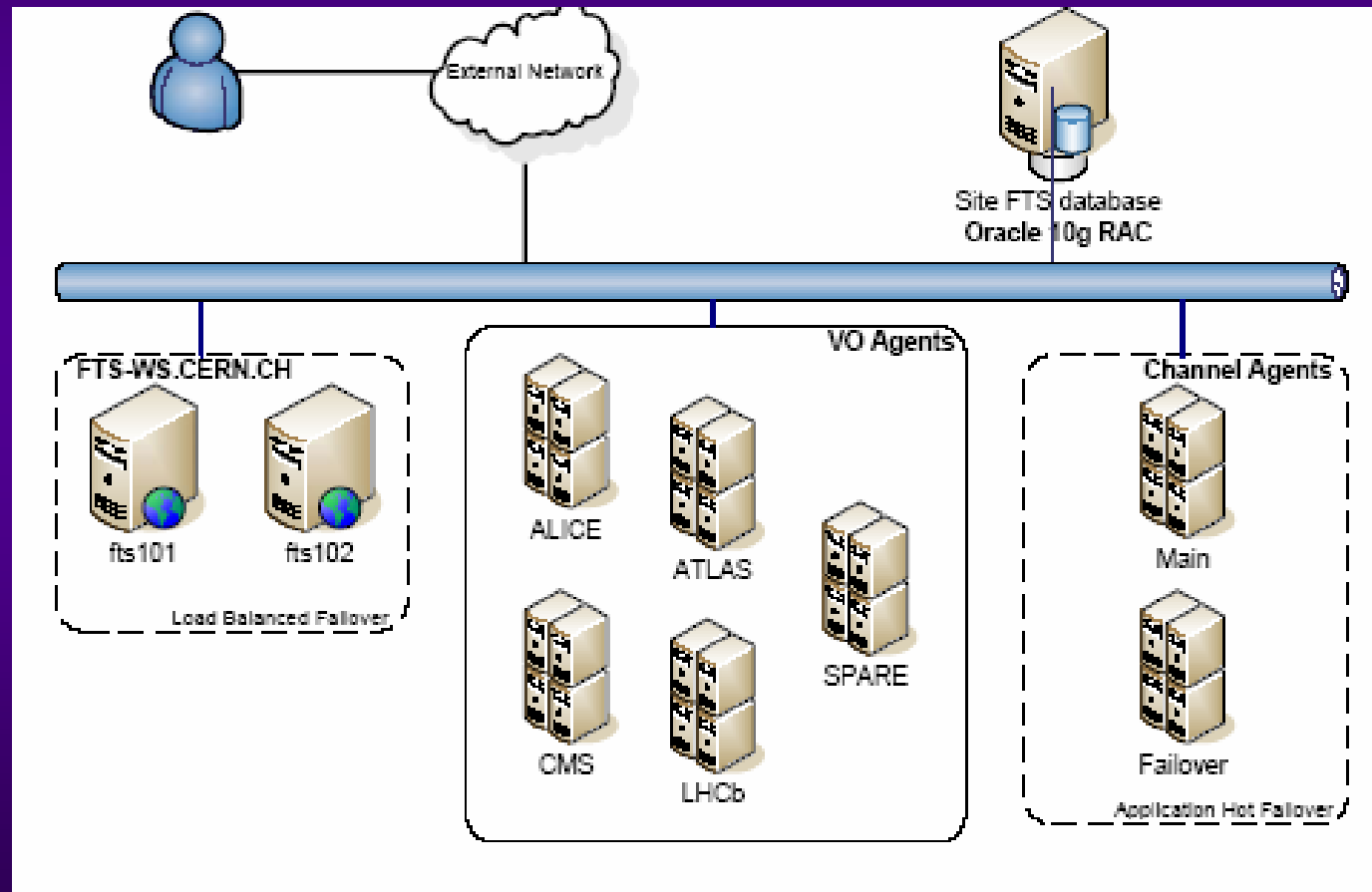
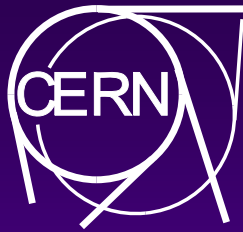
SPARE

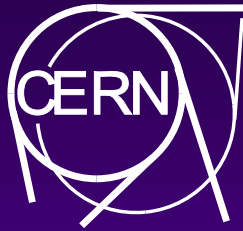


CMS



LHCb

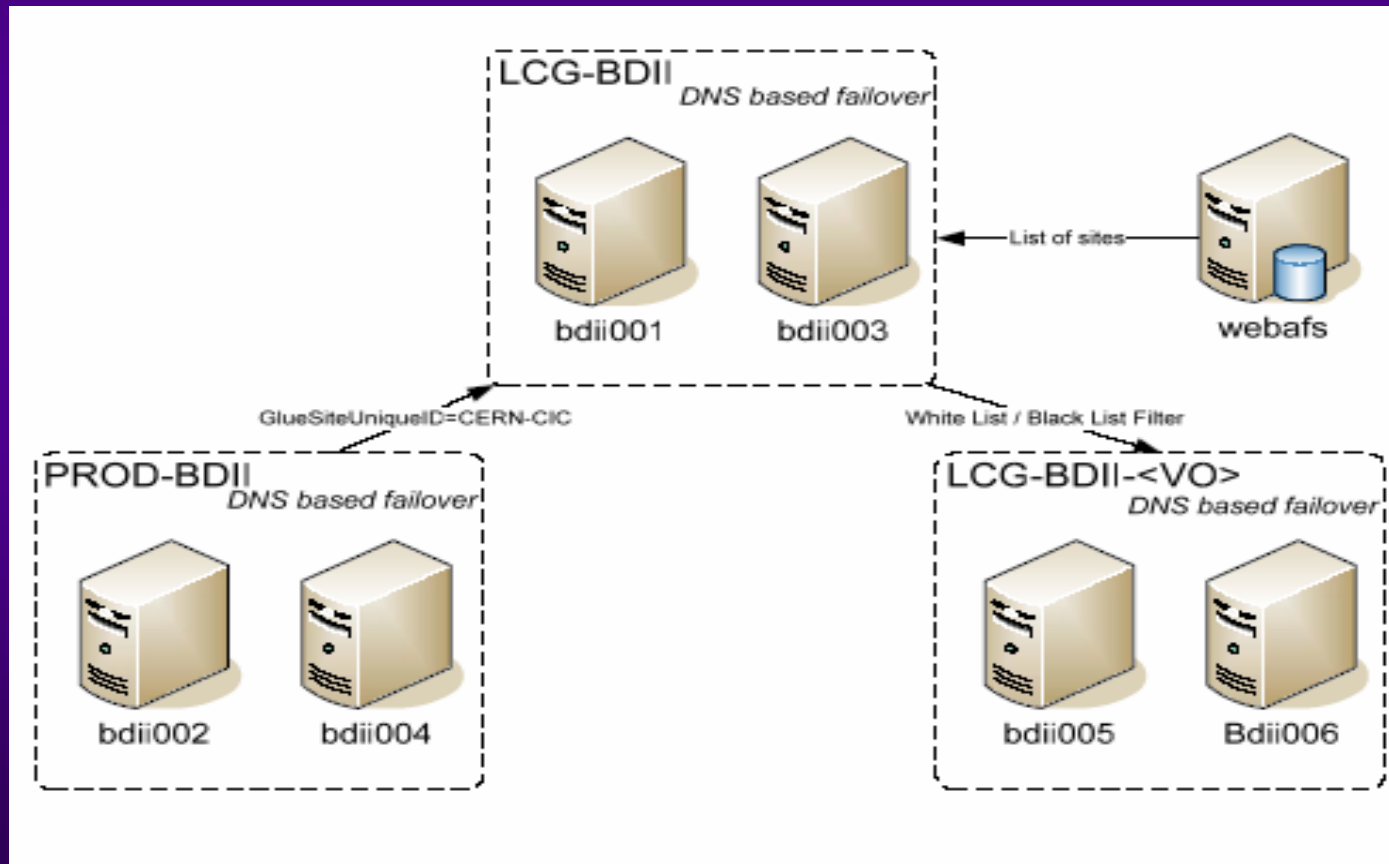
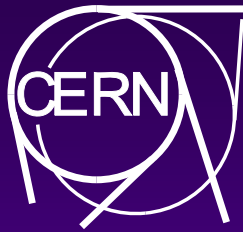


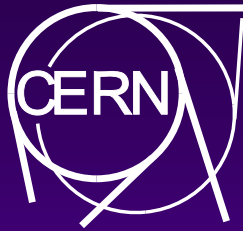


IS: BDII

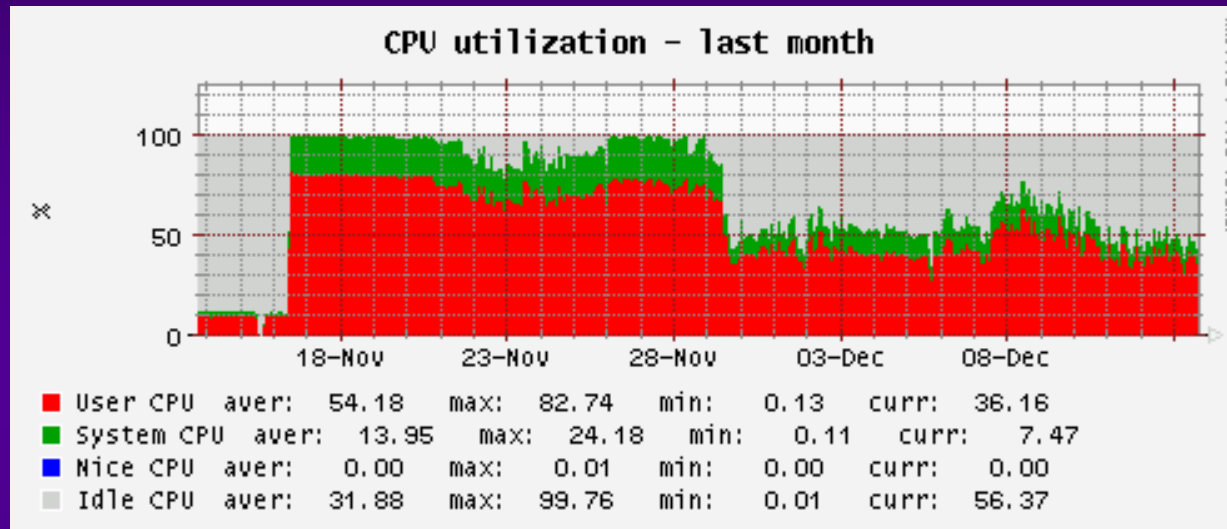
BDII: Load balanced Mid-Range-Servers:

- LCG-BDII top level BDII
- PROD-BDII site BDII
- In preparation: <VO>-BDII
- State information in 'volatile' cache re-queried within minutes

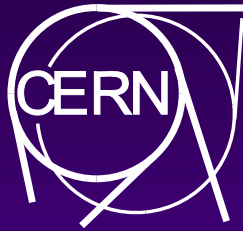




Example: BDII

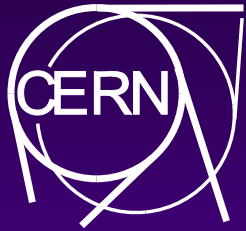


- 17.11.2005 -> First BDII in production
- 29.11.2005 -> Second BDII in production

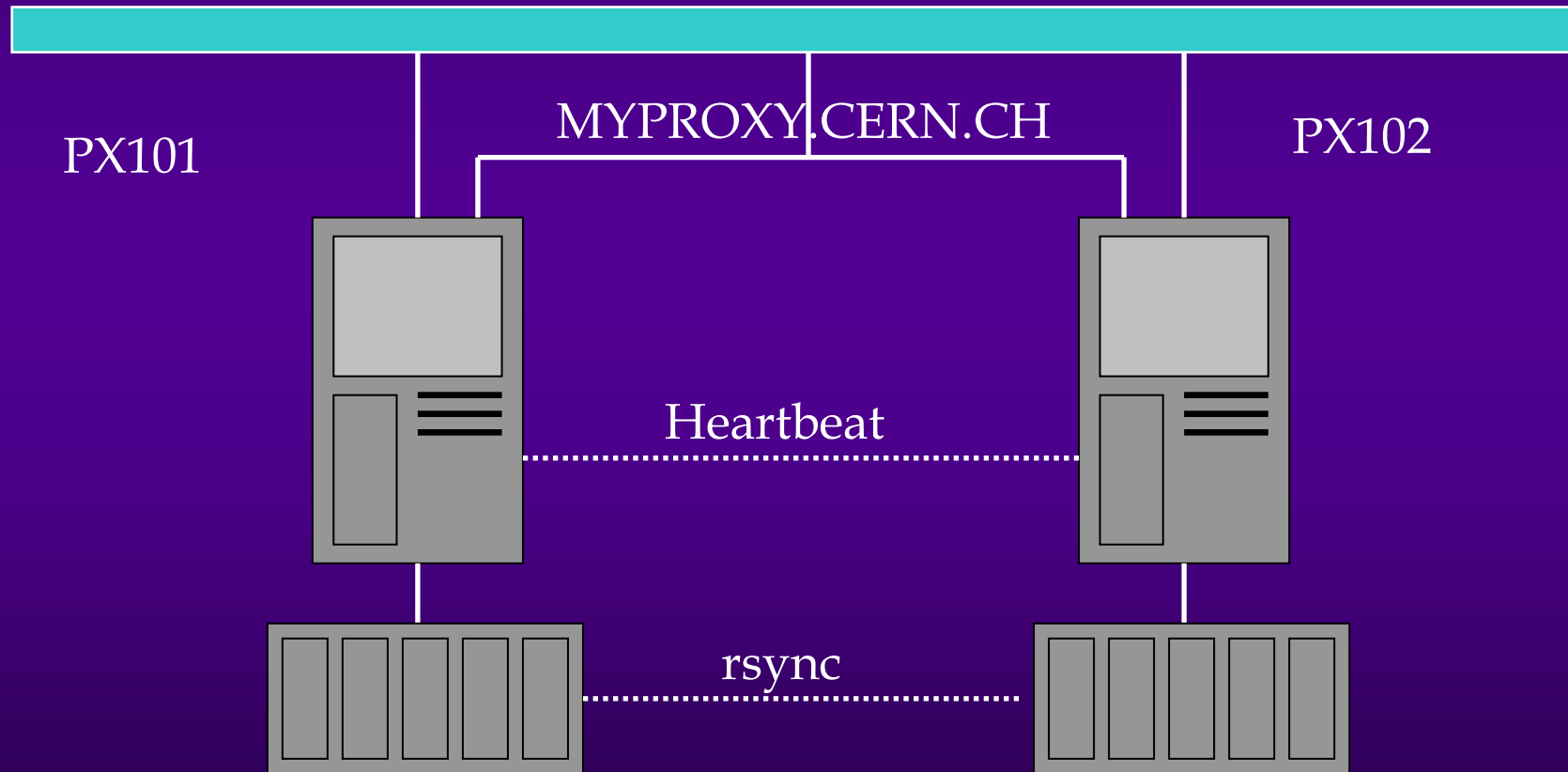


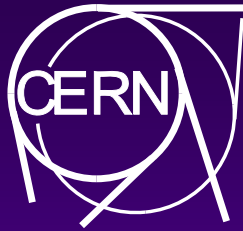
AAS:MyPROxy

- MyProxy has a replication function for a Slave Server, that allows read-only proxy retrieval [not used any more]
- Slave Server gets 'read-write' copy from Master regularly to allow DNS switch over (rsync, every minute)
- HALinux handles failover

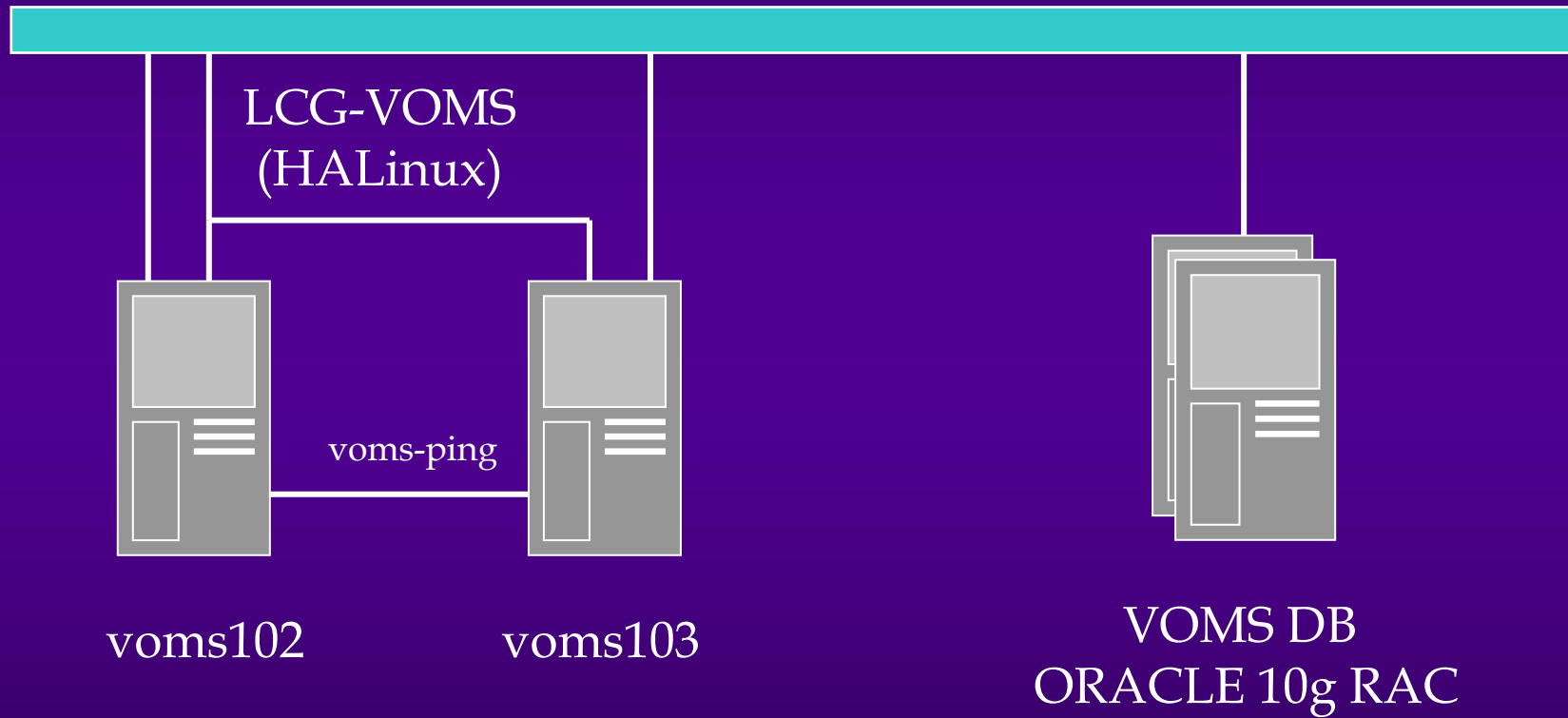


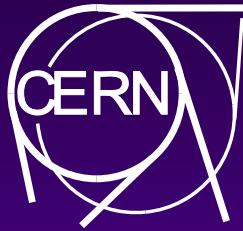
HA Linux





AAS: VOMS





MS: SFT, GRVW, MONB

- Not seen as critical
- Nevertheless servers are migrated to 'Mid-Range-Servers'
- GRVW has a 'hot-standby'
- SFT & MONB on Mid-Range-Servers