



Deploying Production GRID Servers & Services at CERN

Thorsten Kleinwort
CERN IT
DESY Zeuthen 16.05.2006

1



Introduction

- How to deploy GRID services
In a production environment
- h/w used for GRID services
- Various techniques &
First experiences &
Recommendations

16 May 2006

Thorsten Kleinwort CERN-IT

2



Production Servers/Services

What does production mean (for LCG)?

- Reach MoU agreements



MoU

- MoU Annex 3.3 for Tier2:



Production Servers/Services

What does production mean (for LCG)?

- Reach MoU agreements
- Avoid single points of failure for the services:
 - In terms of hardware
 - In terms of software
 - In terms of people
- 'Procedurise':
 - Production/Failure recover/maintenance/...

16 May 2006

Thorsten Kleinwort CERN-IT

5



General remarks

- Production Service Managers † Software/Service Developer/Guru
- Procedure needed for:
S/w installation, upgrade, checking,...
- Requirements needed for:
h/w: disk, memory, backup, high availability
- Operations: 24/7, 9-5, piquet, problem escalation

16 May 2006

Thorsten Kleinwort CERN-IT

6



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

16 May 2006

Thorsten Kleinwort CERN-IT

7



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 May 2006

Thorsten Kleinwort CERN-IT

8



16 May 2006

Thorsten Kleinwort CERN-IT

9



Various Techniques

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

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

16 May 2006

Thorsten Kleinwort CERN-IT

10



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
 - 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)

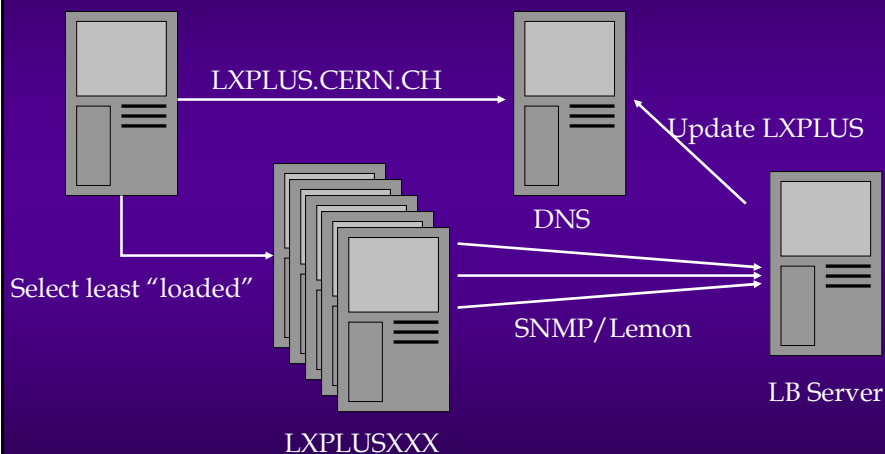
16 May 2006

Thorsten Kleinwort CERN-IT

11



DNS Load Balancing



16 May 2006

Thorsten Kleinwort CERN-IT

12



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 for 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

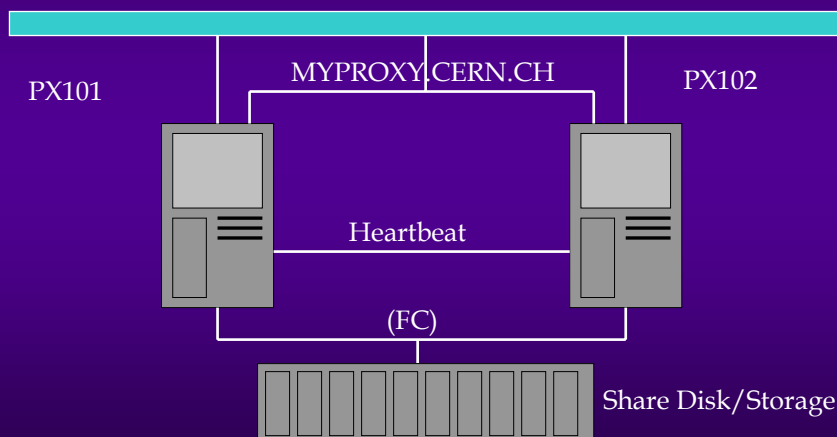
16 May 2006

Thorsten Kleinwort CERN-IT

15



HA Linux



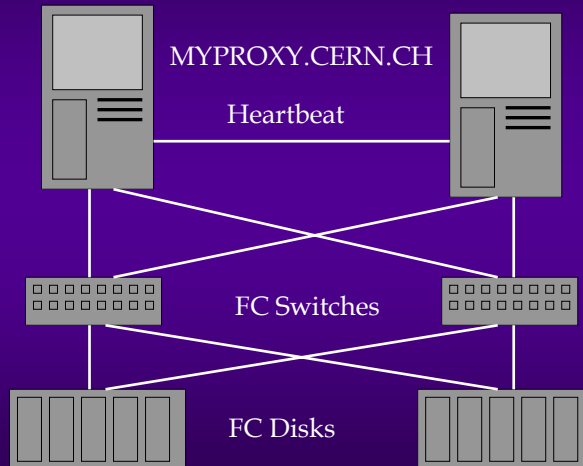
16 May 2006

Thorsten Kleinwort CERN-IT

16



HA Linux + FC disk



No single
Point of
failure

16 May 2006

Thorsten Kleinwort CERN-IT

17



GRID Services: WMS

WMS (Workload Management System):

- WN
- UI
- CE
- RB
- (GRPK)

16 May 2006

Thorsten Kleinwort CERN-IT

18



WMS: WN & UI

WN & UI:

- Simple, add as many machines as you need
- Do not depend on single point of failures
- Service can be kept up while interventions are ongoing:
Have procedures for s/w, kernel, OS upgrades



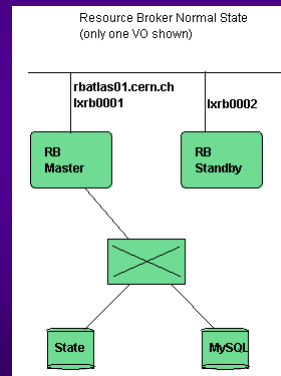
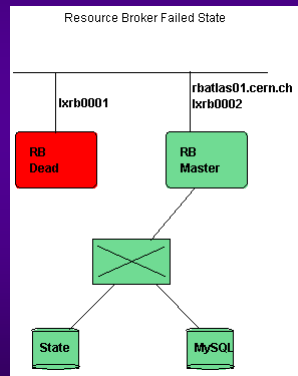
WMS: CE & RB

CE & RB are stateful Servers:

- load balancing only for scaling:
 - Done for RB and CE
- CE's on 'Mid-Range-Server' h/w
- RB's now done with gLite 3.0
- If one server goes, the information it keeps is lost
- Possibly FC storage backend behind load balanced, stateless servers



Failover with FC



16 May 2006

Thorsten Kleinwort CERN-IT

21



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

16 May 2006

Thorsten Kleinwort CERN-IT

22



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

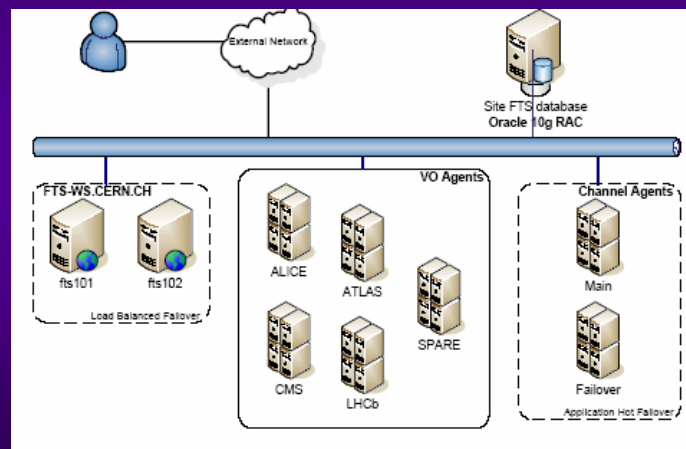


LHCb

16 May 2006

Thorsten Kleinwort CERN-IT

25



16 May 2006

Thorsten Kleinwort CERN-IT

26



IS: BDII

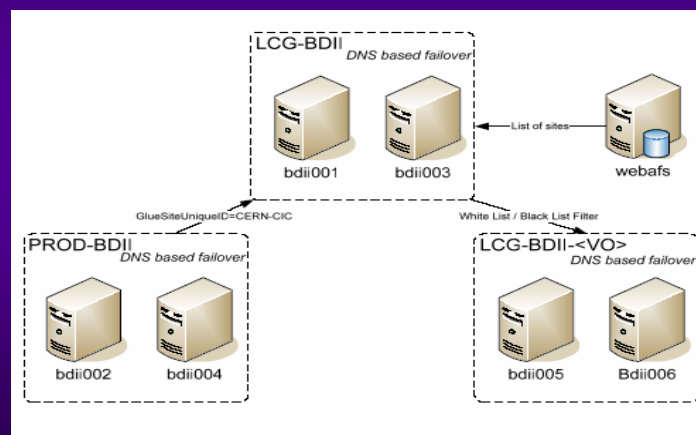
BDII: Load balanced Mid-Range-Servers:

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

16 May 2006

Thorsten Kleinwort CERN-IT

27



16 May 2006

Thorsten Kleinwort CERN-IT

28



Grid Services: AAS & MS

AAS: Authentication & Authorization Services

- PX: (My)Proxy Service
- VOMS: Virtual Organisation Membership Service

MS: Monitoring Service

- MONB: RGMA
- GRVW: GridView
- SFT: Site Functional Tests

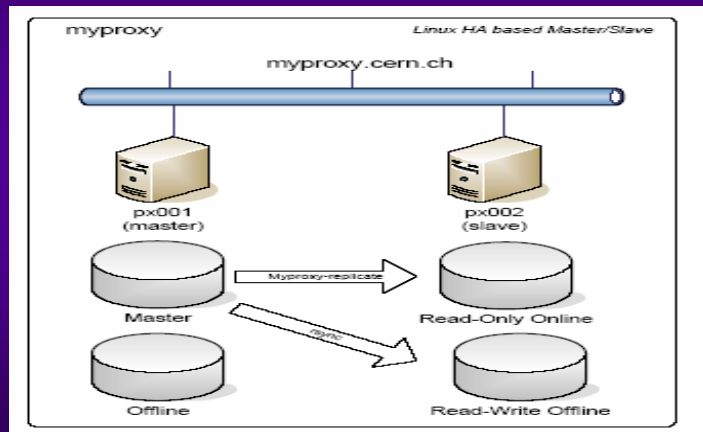
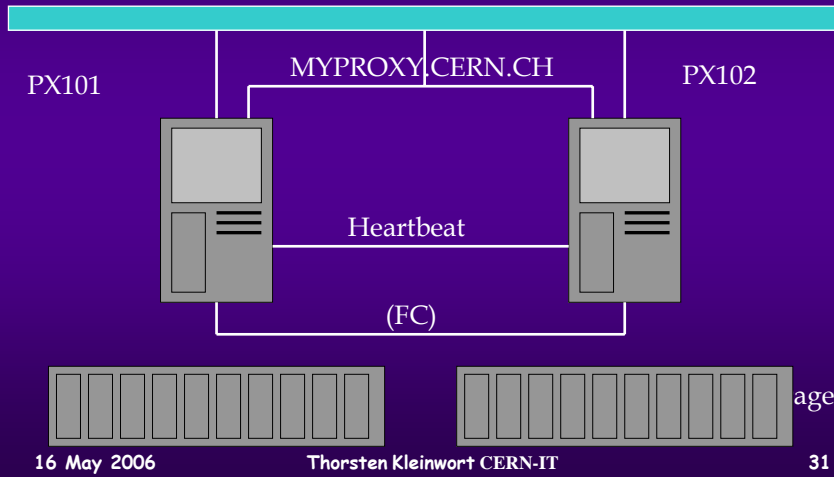


AAS: MyPROxy

- MyProxy has a replication function for a Slave Server, that allows read-only proxy retrieval
- Slave Server gets 'read-write' copy from Master regularly to allow DNS switch over
- HALinux handles failover



HA Linux



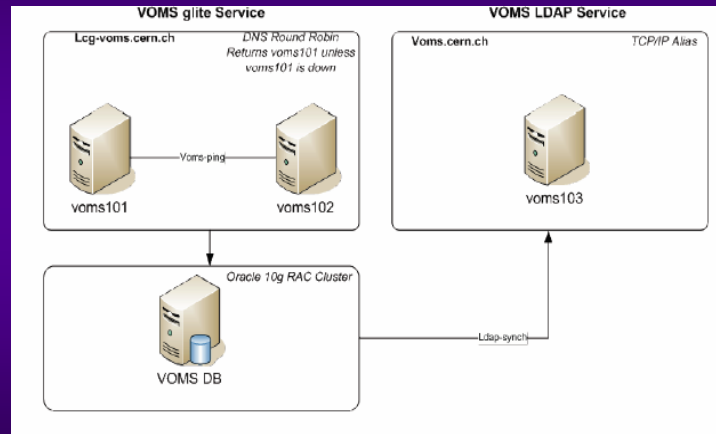
16 May 2006

Thorsten Kleinwort CERN-IT

32



AAS: VOMS



16 May 2006

Thorsten Kleinwort CERN-IT

33



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

16 May 2006

Thorsten Kleinwort CERN-IT

34



Other Experiences

- Sometimes difficult to integrate the GRID middleware with the local fabric
 - RPM dependencies
 - Monitoring
 - Network connectivity
 - UI incompatible
- Scalability issues
 - CE already overloaded

16 May 2006

Thorsten Kleinwort CERN-IT

35



Monitoring

- (Local) Fabric Monitoring (Lemon)
- GRID Server Monitoring
- GRID Service Monitoring
- GRID Monitoring:
Gstat, SFT, GRIDVIEW

16 May 2006

Thorsten Kleinwort CERN-IT

36



People

- Operators -> SysAdmins -> Service Managers (Fabric+Service) -> Application Managers
- Regular daily Meeting (Merge of two)
- Regular weekly meeting (CCSR)
- WLCG weekly conference call
- Procedures

16 May 2006

Thorsten Kleinwort CERN-IT

37



Conclusions

- We are on track to have the current GRID services 'production quality' for SC4
- As much as we can do (and want to effort) on the Fabric Level:
h/w, s/w, people
- Still 'wish list' for improvements of the services

16 May 2006

Thorsten Kleinwort CERN-IT

38