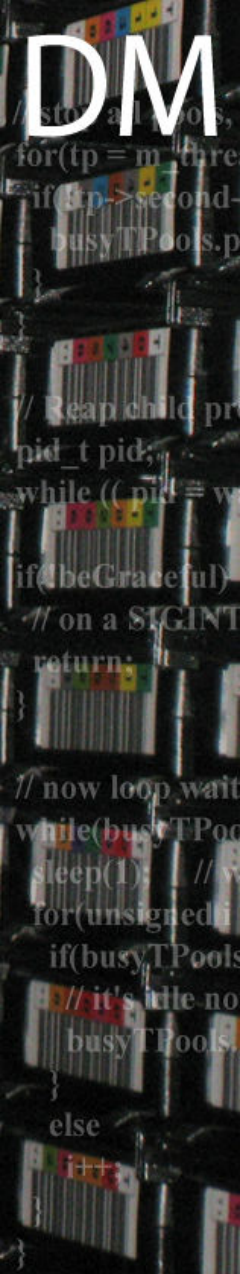


Oracle 11g R2 **New** Features RAC and ASM



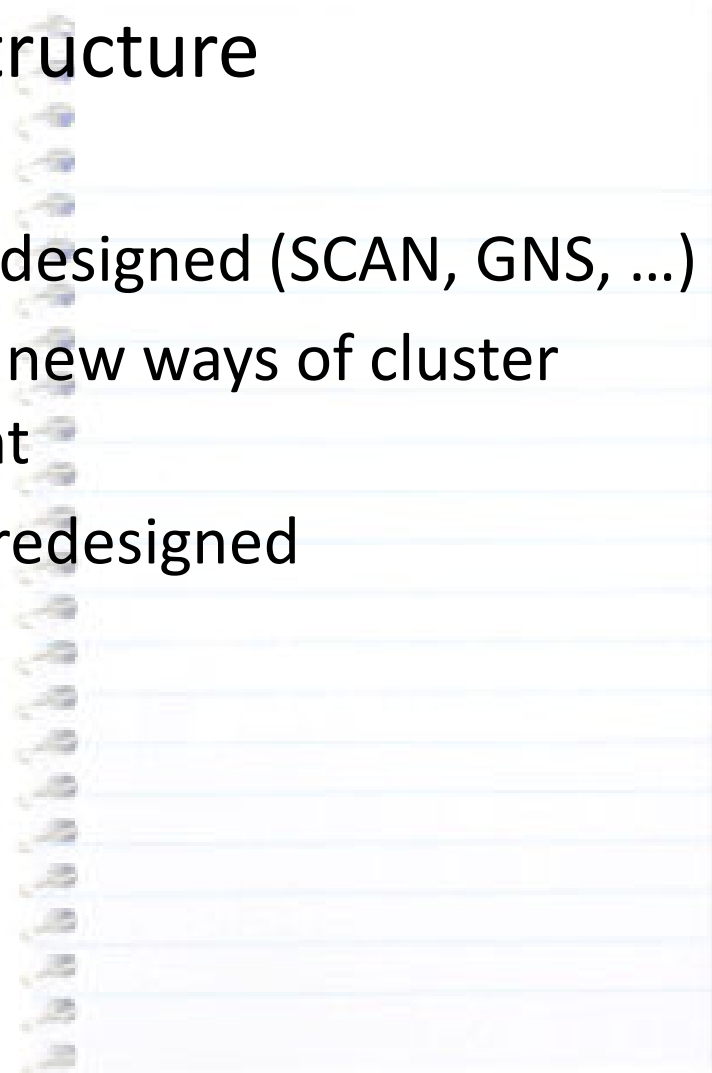
Dawid Wojcik

27 November 2009



DM Outline

- The **new** Grid Infrastructure
 - Introduction
 - Oracle connectivity redesigned (SCAN, GNS, ...)
 - Cluster resources and new ways of cluster database management
 - OCR and voting disks redesigned
- ACFS and ADVMM
 - Features
 - Architecture
 - Usage scenarios
- Fun stuff



- Clusterware is now called *Oracle Grid Infrastructure*



- **Installation**

- ASM is now part of Grid Infrastructure
- OCR and voting disk no longer supported on raw or block devices, only on ASM
- If installation goes wrong Oracle now provides a **deinstallation utility**

- **Upgrade**

- Rolling upgrade supported only from 11.1 (for 10g DB is required to be down on all nodes)
- Possible to install 11.2 Grid with 10.2 DBs (you need to “pin” the nodes)

- Grid Infrastructure has been redesigned
 - New process tree
 - *OPROCD*, *OCLSOMON*, *OCLSVMON* no longer exist
 - Init spawns only one process (*ohasd*), which spawns
 - oraagent
 - orarootagent
 - cssdagent
 - Hangcheck timer no longer required on linux
 - ASM instance is started during Grid Infrastructure startup to allow accessing OCR and voting disks (*bootstrap* procedure)



- **SCAN** (Single Client Access Name)
 - Allows any client to connect to any database or service in the cluster with just **single hostname** even if the cluster size changes
 - **Load balances** connections across all instances providing a service (server side load balancing)
 - **Provides failover** in case an instance fails and services relocate
 - Allows clients to use EZConnect or simple JDBC connection in 'RAC-aware' mode
 - Each cluster (even two node one) has
 - 3 SCAN VIPs
 - 3 SCAN listeners



- **SCAN** implementation



- **DNS** (Domain Name Service)

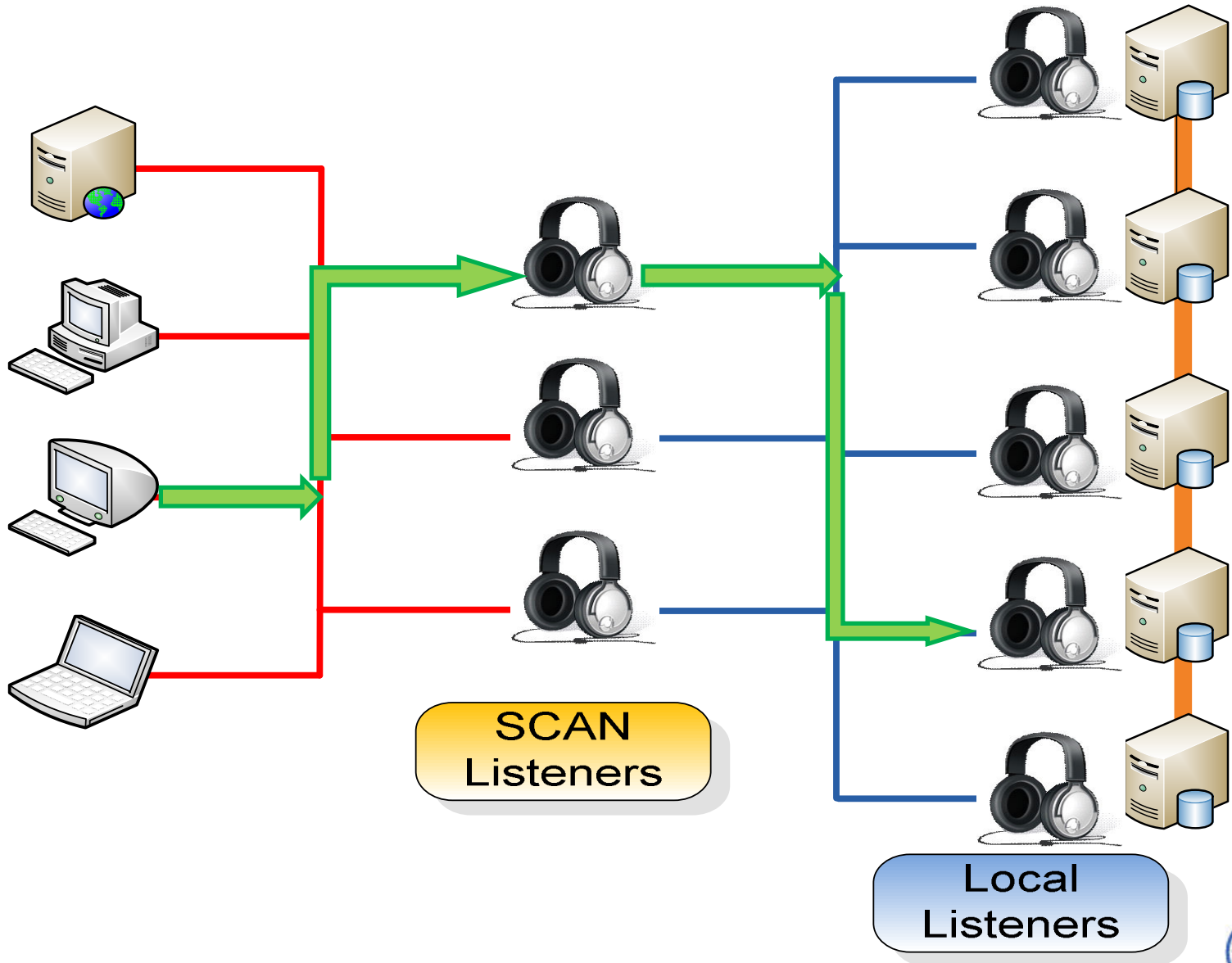
- NetOps define DNS alias which resolves one of three SCAN VIPs in round-robin way


- **GNS** (Grid Naming Service)

- Subdomain DNS resolution for SCAN VIPs
- Can be used for DHCP VIP registration

- **SCAN registration**

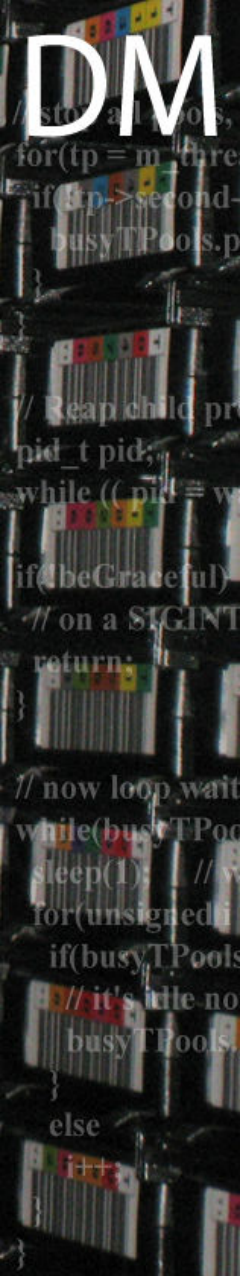
- Instance **registers** with **local listener** on its node
- Instance **registers** with SCAN (all **SCAN listeners**)



- **Administrator Managed** – “the old way” 
 - DBA defines where the database should run and manages list of nodes in RAC
 - DBA defines where services should run within the cluster
- **Policy Managed** – “the new way”
 - DBA defines resource requirements for each DB running in the cluster
 - DB instances are started “when needed”
 - Goal - remove hard coding services and instance nodes

- Allows **grouping** more **DBs** into one cluster for better **resource utilization**
 - Cluster divided into **pools of servers**
 - Applications and DBs run in one or more server pools
 - Instances can be started on different nodes each time the cluster starts (!)
 - **Resource allocation** is defined by 3 **attributes**:
 - **Min** – minimum number of servers (default 0)
 - **Max** – maximum number of servers (default 0 or -1)
 - **Importance** – 0 (least important) to 1000

- Services in RAC managed with server pools can run only in one server pool
- Depending on how many nodes they run they can be
 - **Uniform** – run in all instances in a pool
 - **Singleton** – run only in one instance in a pool
- This gives less flexibility than Administrator Managed Services

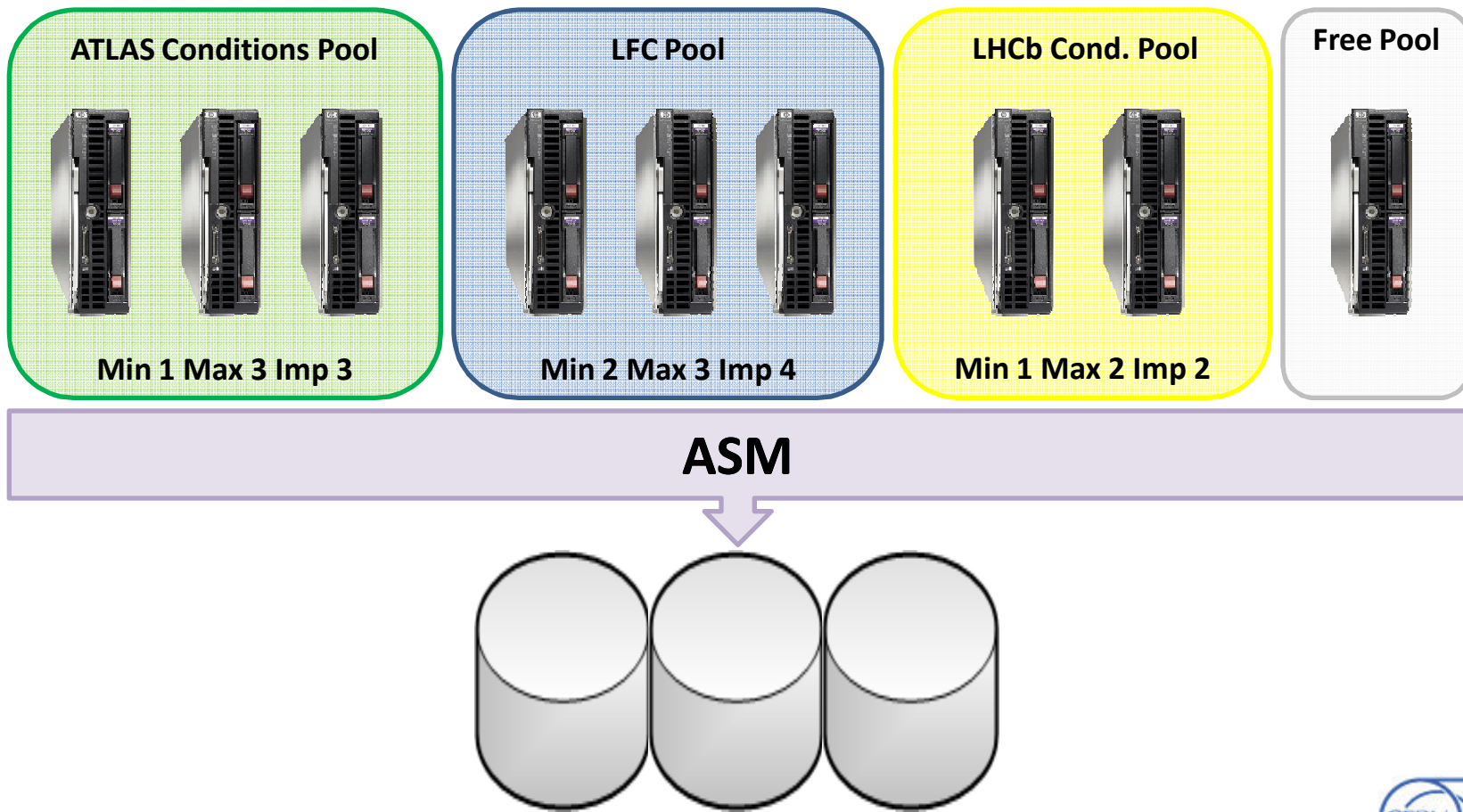


Server assignment in Grid

- Servers are assigned in the following **order**:
 - **Generic server pool** (if defined, during upgrade all servers go to generic pool)
 - **User defined server pools**
 - **Free pool**
- Grid Infrastructure uses **importance** of server pool to determine server assignment order:
 - Fill all server pools in order of importance until they meet their **minimum**
 - Fill all server pools in order of importance until they meet their **maximum**
 - By Default any left over go into **free pool**

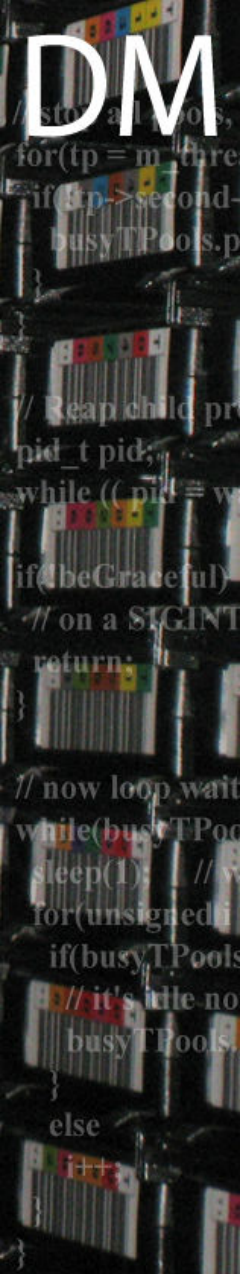


- Cluster of 9 nodes is starting up ...



- If a server **leaves** the cluster ... and the free pool is empty
 - Grid Infrastructure may **move servers** from one server pool to another **only if**
 - you have **non-default** values for **min**, **importance**
 - a server **pool falls below its minimum**
 - Servers to be **moved** are taken from
 - A server pool that is **less important**
 - A server pool of the same importance which has **more servers than its minimum**

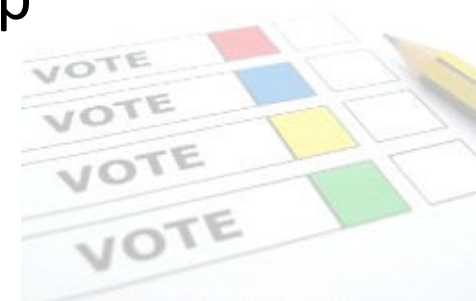




- **OCR** is now just like any **file in ASM**
 - New Oracle managed file type
 - **1 OCR per diskgroup**
 - Good practice to have more than 1 OCR
 - OCR file redundancy follows diskgroup redundancy
 - **OCR now holds automatic backups of voting disks!**
(taken automatically on any cluster configuration change)



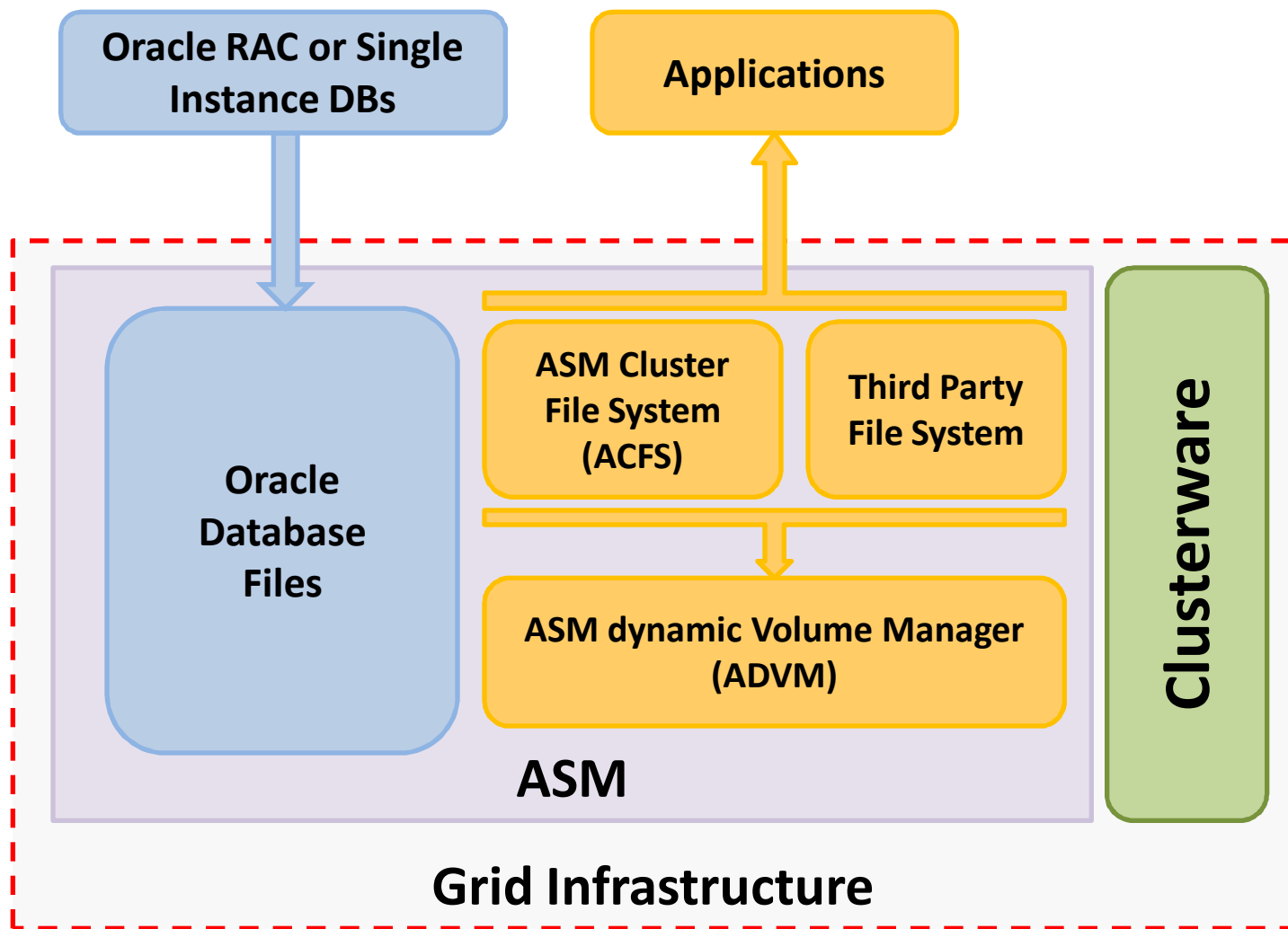
- **Voting Disks** are created on specific disks and Grid Infrastructure knows their location
 - Number of voting disks depends on the redundancy chosen for the diskgroup
 - External – 1 Voting Disk
 - Normal – 3 Voting Disks
 - High – 5 Voting Disks
 - Diskgroup must contain enough failure groups to create each voting disk in a separate failure group!

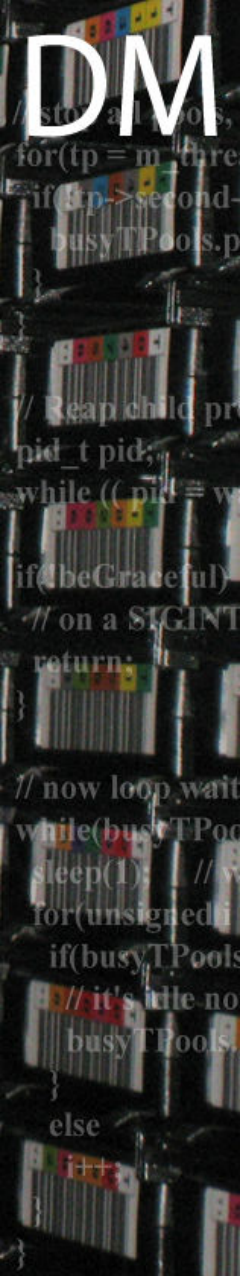




- What is ACFS?

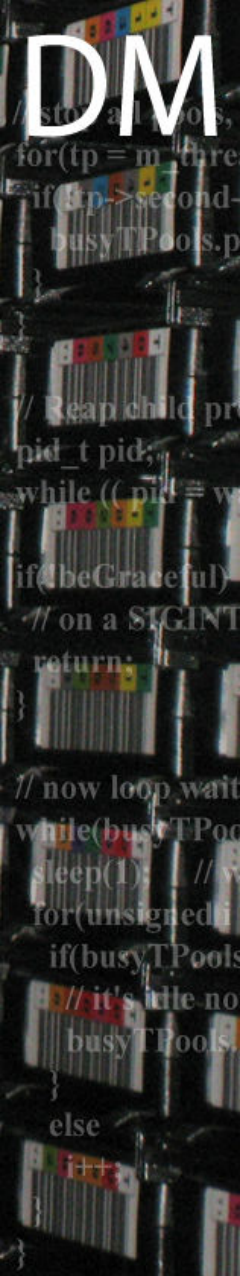
- ASM Cluster File System – POSIX/X-OPEN compliant File System
 - Can be used cluster-wide or single node only
- Currently Linux only (RHEL5 +), new platforms coming soon
- Can be shared using NFS, CIFS, ...
- **Online** filesystem **expansion / shrink**
- mirror protection when using NORMAL redundancy diskgroups
- **read-only snapshots built-in**





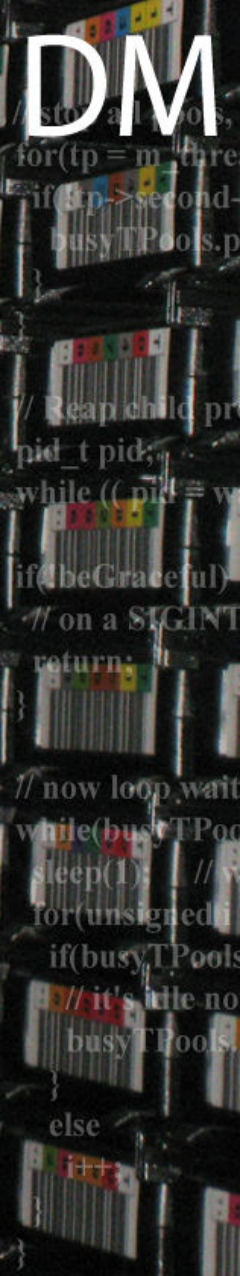
- ASM Dynamic Volume Manager (**ADVM**)
 - Carves out logical volumes from an **ASM diskgroup** that will be exposed to **OS** as **block devices**
 - Device name format: `/dev/asm/volume_name-xyz`
 - ADVM volume can be partitioned
 - On top of ADVM volume one can create any file system (ext3, ACFS, ...)
 - Volumes can be resized (make sure the filesystem supports it)
- Kernel modules: `oracleacfs`, `oracleadv`, `oracleoks`, `oracleasm`

- ADVM Volumes can be exposed (enabled) on any nodes
- **Volume redundancy** can be
 - **Unprotected** – for external redundancy diskgroup
 - **High, mirror or unprotected** – for normal redundancy diskgroup
- Stripe width and IDP (Intelligent Data Placement) parameter can be specified



ASM – some new features

- **IDP (Intelligent Data Placement)**
 - Requires JBOD storage connected (one disk – one LUN)
 - ASM files can be specified to be on external (hot, faster) part of the disks or internal (cold, slower)
- New tools to manage ASM and ACFS
 - Extended asmcmd tool
 - New tool for ACFS – acfsutil
- ...

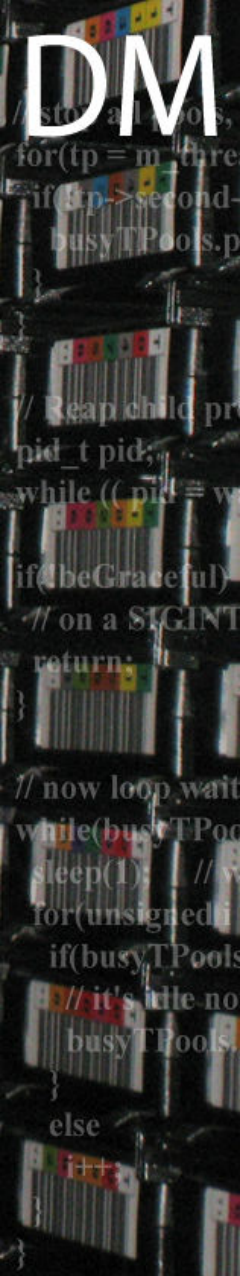


ACFS use cases

- BFILEs, external tables, external files storage
- Local file systems, exports, dumps, etc.
- Batch jobs data
- Shared RDBMS home (not recommended if using rolling patches)
- Local RDBMS home
- ...

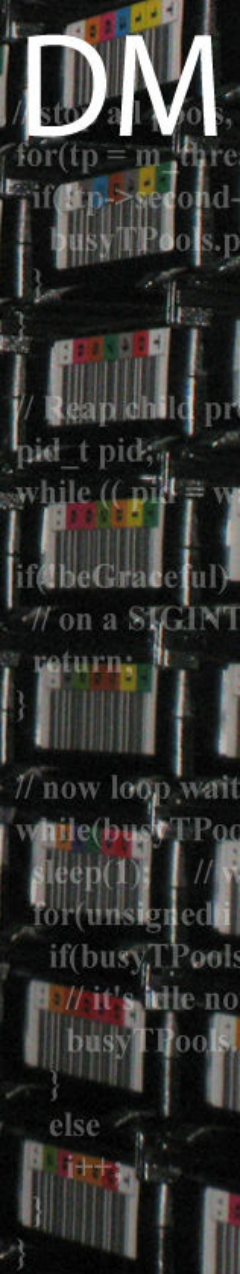


- No need for 3rd party cluster file systems
- Uses ASM - known tools
- it will become platform independent CFS
- Simple administration (acfsutil, acfschkdsk, mkfs, ...)
- Cost effective (if you have at least one node licensed for SE or EE - you can use it)
- Snapshots for free
- Very good performance



ACFS snapshots

- Provide Point-In-Time images (**ACFS only**)
- **File System functionality**
- Can be performed online
- Used for consistent backups
- Limited to 63 snapshots per File System
- **Copy on write mechanism** (before-images shared between snapshots)
- Snapshots within the same file system - monitor free space!



- Compared ext3 on local disk vs. ext3 on ADVM
 - No difference in write speed
- Typical read/write performance is **much better** comparing to **ext3** (both running on ADVM)

Test type	EXT3 (on ADVM)	ACFS (on ADVM)
Big file creation	37MB/s	225MB/s
Extracting 2,5GB tar	180 s	100 s
Deleting big dir. tree	3 s	11 s
Touching 50k files	63 s	75 s



- Recursive Subquery Factoring

```
with x( s, ind ) as
( select sud, instr( s, '537' ) as ind
  from ( select '53798648476624195879'
        sud from dual )
  union all
  select substr( s, 1, ind-1 ) ||
        , instr( s, '537' ) as ind
  from x
  , ( select to_char( rownum, '00' )
      from dual
      connect by rownum <= 9 ) z
  where ind > 0
  and not exists ( select 1
                  from x
                  where s = substr( s, 1, ind-1 ) || z
                  or s = substr( s, 1, ind+1 ) || z
                  or s = substr( s, ind+1, ind-1 ) || z
                )
)
```

5	3			7				
6			1	9	5			
	9	8					6	
8				6				3
4			8		3			1
7				2				6
	6					2	8	
			4	1	9			5
				8			7	9

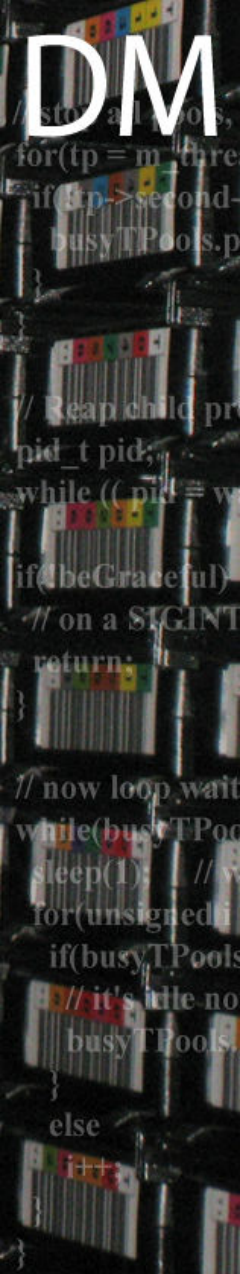
28 419 5 8 79'

1)
* 3

<http://technology.amis.nl/blog/6404/oracle-rdbms-11gr2-solving-a-sudoku-using-recursive-subquery-factoring>

```
)
select s
from x
where ind = 0
/
```





DM Conclusions

- 11gR2 comes with many new features
 - We will test some of the most interesting and publish tests on our TWiki
 - We are eager to try ACFS in production
 - Not ready for Server Pool managed clusters (isolation preferred over consolidation)
 - Except for downstream capture boxes
 - ... more useful features to be discovered