# 10g RAC Best Practices

Kirk McGowan
Technical Director – RAC Pack

Server Technologies

Oracle Corporation

**ORACLE**®

# Disclaimer

These Best Practices are based on customer experiences, and they will generally give the best results. However, systems have different requirements and cost structures, so these Best Practices might not be applicable in all cases. As technology evolves and with new experiences, these Best Practices will probably change over time. These Best Practices do not replace the standard product documentation which is the official guide to product use.
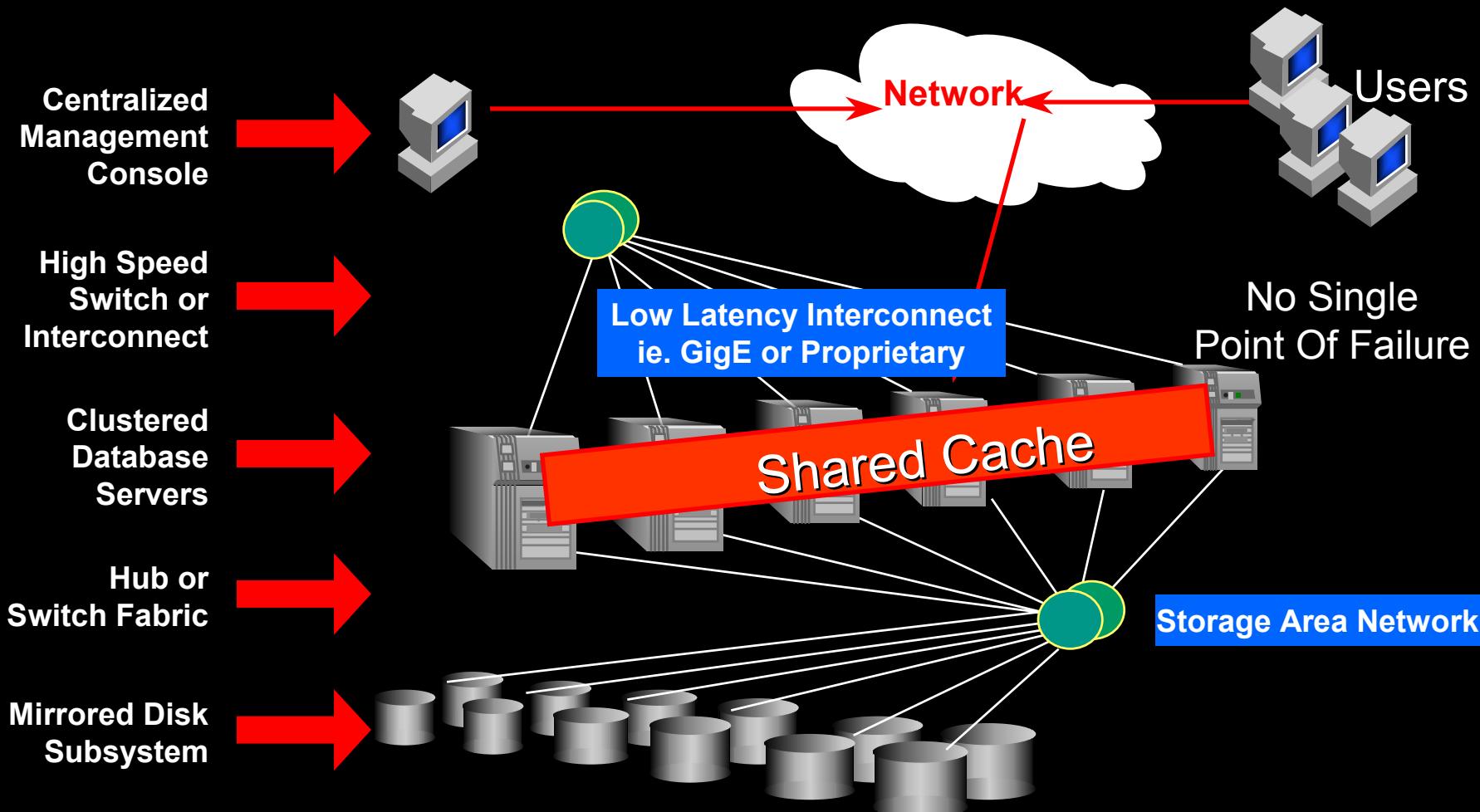
# Agenda

- Planning Best Practices
  - Understand and Plan the Architecture
  - Manage Expectations
  - Define objectives and success criteria
  - Project plan

- Implementation Best Practices
  - Infrastructure considerations
  - Installation/configuration
  - Database creation
  - Application considerations

- Operational Best Practices
  - Backup & Recovery
  - Performance Monitoring and Tuning
  - Production Migration

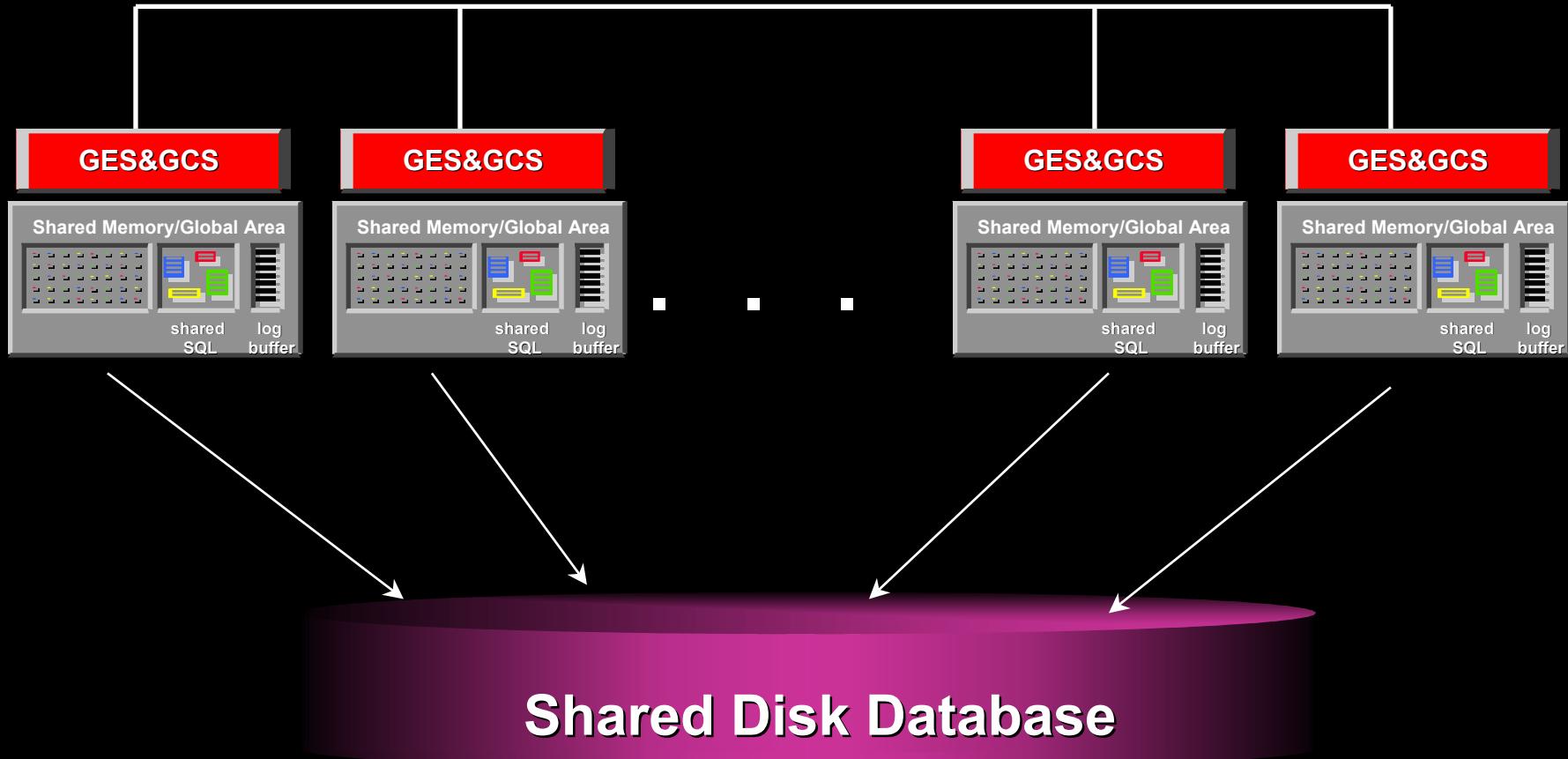**ORACLE**

# Planning

- Understand the Architecture
  - Cluster terminology
  - Functional basics
    - HA by eliminating node & Oracle as SPOFs
    - Scalability by making additional processing capacity available incrementally
  - Hardware components
    - Private interconnect/network switch
    - Shared storage/concurrent access/storage switch
  - Software components
    - OS, Cluster Manager, DBMS/RAC, Application
    - Differences between cluster managers

**ORACLE**

# RAC Hardware Architecture

**Centralized Management Console** →

**High Speed Switch or Interconnect** →

**Clustered Database Servers** →

**Hub or Switch Fabric** →

**Mirrored Disk Subsystem** →

**Network**

Users

No Single Point Of Failure

**Low Latency Interconnect ie. GigE or Proprietary**

**Shared Cache**

**Storage Area Network**

**ORACLE**

# RAC Software Architecture

## Shared Data Model



GES&GCS — Shared Memory/Global Area — shared SQL — log buffer

GES&GCS — Shared Memory/Global Area — shared SQL — log buffer

GES&GCS — Shared Memory/Global Area — shared SQL — log buffer

GES&GCS — Shared Memory/Global Area — shared SQL — log buffer

**Shared Disk Database**

# 10g Technology Architecture

# Plan the Architecture

- Eliminate SPOFs
  - Cluster interconnect redundancy (NIC bonding/teaming, …)
  - Implement multiple access paths to the storage array using 2 or more HBA's or initiators
    - Investigate multi-pathing sw over these multiple devices to provide load balancing and failover.
- Processing nodes – sufficient CPU to accommodate failure
- Scalable I/O Subsystem
  - Scalable as you add nodes
- Workload Distribution (load balancing) strategy
  - Net Services (SQL*Net)
  - Oracle10g Services
- Establish management infrastructure to manage to Service Level Agreements
  - Grid Control

**ORACLE**

# Cluster Hardware Considerations

- Cluster interconnects
  - FastEthernet, Gigabit Ethernet, Proprietary interconnects (SCI, Hyperfabric, memory channel, …)
  - Dual interconnects, stick with GigE/UDP
- Public networks
  - Ethernet, FastEthernet, Gigabit Ethernet
- Server Recommendations
  - Minimum 2 CPUs per server
  - 2 and 4 CPU servers normally most cost effective
  - 1-2 GB of memory per CPU
  - Dual IO Paths
- Intelligent storage, or JBOD
- Fiber Channel, SCSI, iSCSI or NAS storage connectivity
- Future: Infiniband

# Plan the Architecture

- Shared storage considerations (ASM, CFS, shared raw devices)
- Use S.A.M.E for shared storage layout
    - http://otn.oracle.com/deploy/availability/pdf/oow2000_same.pdf
- Local ORACLE_HOME versus shared ORACLE_HOME
- Separate HOMEs for CRS, ASM, RDBMS
- OCR and Voting Disk on raw devices
    - Unless using CFS

**ORACLE**

# RAC Technology Certification

- For more details on software certification and compatible hardware:

  – http://technet.oracle.com/support/metalink/content.html

- Discuss hardware configuration with your HW vendor

- Try to stick to standard components that have been properly tested/certified

ORACLE

# Set Expectations Appropriately

- If your application will scale transparently on SMP, then it is realistic to expect it to scale well on RAC, without having to make any changes to the application code.

- RAC eliminates the database instance, and the node itself, as a single point of failure, and ensures database integrity in the case of such failures

**ORACLE**

# Planning: Define Objectives

- Objectives need to be quantified/measurable
  - HA objectives
    - Planned vs. unplanned
    - Technology failures vs. site failures vs. human errors
  - Scalability Objectives
    - Speedup vs. scaleup
    - Response time, throughput, other measurements
  - Server/Consolidation Objectives
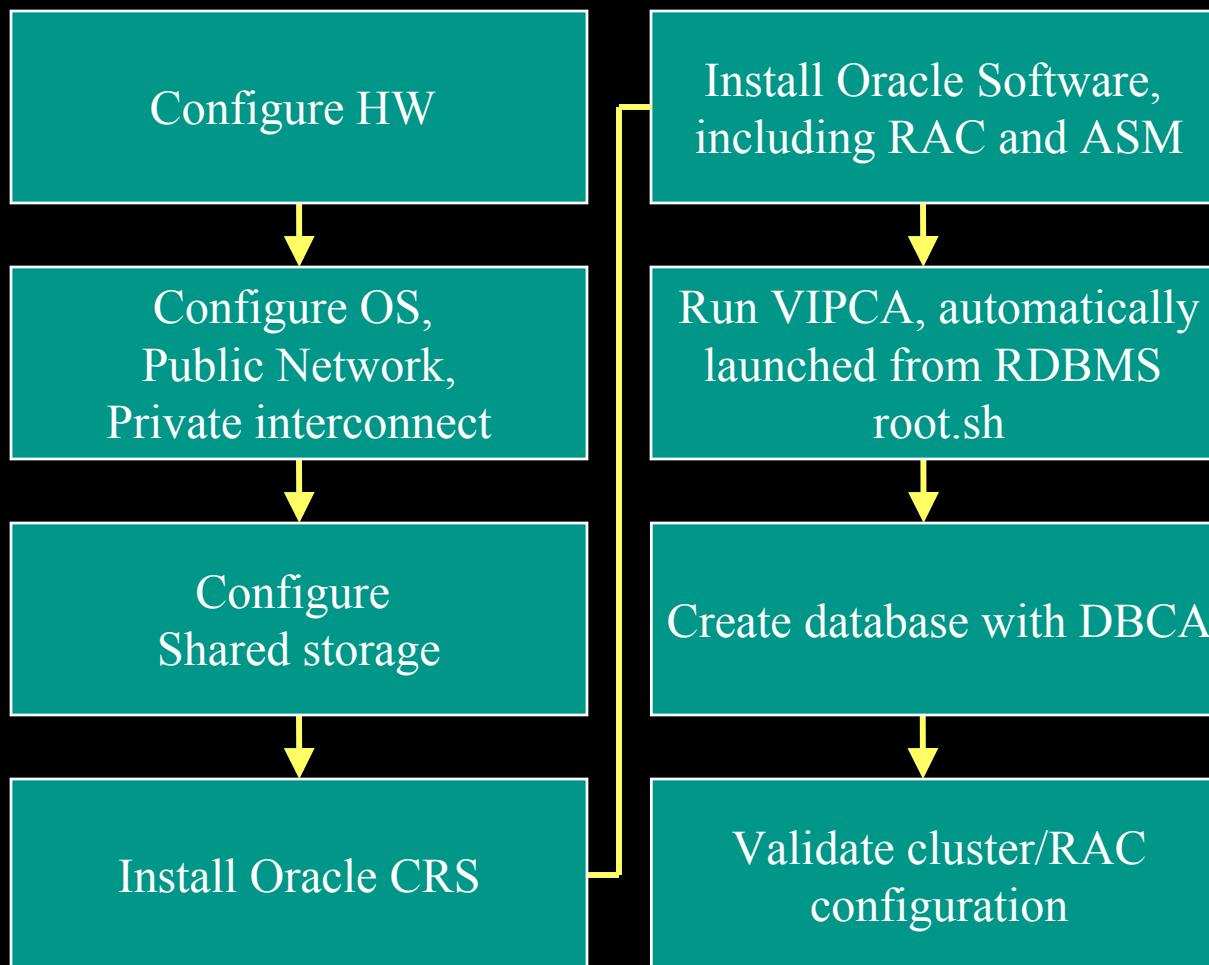    - Often tied to TCO
    - Often subjective

# Build your Project Plan

- Partner with your vendors
  - Multiple stakeholders, shared success
- Build detailed test plans
  - Confirm application scalability on SMP before going to RAC ➔ optimize first for single instance
- Address knowledge gaps and training
  - Clusters, RAC, HA, Scalability, systems management
  - Leverage external resources as required
- Establish strict System and Application Change control
  - Apply changes to one system element at a time
  - Apply changes first to test environment
  - Monitor impact of application changes on underlying system components
- Define Support mechanisms and escalation procedures
  - Including dedicated, long term, test cluster

ORACLE

# Agenda

- Planning Best Practices
    - Architecture
    - Expectation setting
    - Objectives and success criteria
    - Project plan
- **Implementation Best Practices**
    - **Installation/configuration**
    - **Database creation**
    - **Application considerations**
- Operational Best Practices
    - Backup & Recovery
    - Performance Monitoring and Tuning
    - Production Migration

**ORACLE**

# Implementation Flowchart

Configure HW

↓

Configure OS,
Public Network,
Private interconnect

↓

Configure
Shared storage

↓

Install Oracle CRS

→

Install Oracle Software,
including RAC and ASM

↓

Run VIPCA, automatically
launched from RDBMS
root.sh

↓

Create database with DBCA

↓

Validate cluster/RAC
configuration

**ORACLE**

# Operating System Configuration

- Confirm OS requirements from
    - Platform-specific install documentation
    - Quick install guides (if available) from Metalink/OTN
    - Release notes
- Follow these steps on EACH node of the cluster
    - Configure ssh
        - 10g OUI uses ssh, not rsh
    - Configure Private Interconnect
        - Use UDP and GigE
        - Non-routable IP addresses (eg 10.0.0.x)
        - Redundant switches as std configuration for ALL cluster sizes.
        - NIC teaming configuration (platform dependant)
    - Configure Public Network
        - VIP and name must be DNS-registered in addition to the standard static IP information
        - Will not be visible until VIPCA install is complete

ORACLE

# NIC Bonding

- Required for private interconnect resiliency.

- Various 3rd party vendor solutions available:

  - Linux

    - NIC bonding in RHEL 3.0 ES [/http://www.kernel.org/pub/linux/kernel/people/marcelo/linux-2.4/Documentation/networking/bonding.txt](/http://www.kernel.org/pub/linux/kernel/people/marcelo/linux-2.4/Documentation/networking/bonding.txt)
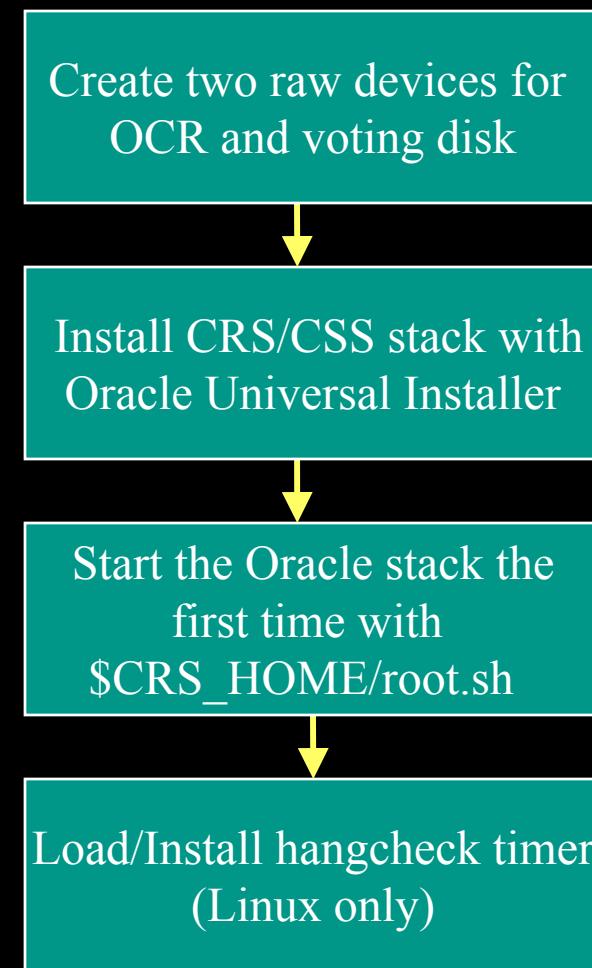
    - Intel® Advanced Network Services (ANS) [http://www.intel.com/support/network/adapter/1000/linux/ans.htm](http://www.intel.com/support/network/adapter/1000/linux/ans.htm)

    - HANIC [http://oss.oracle.com/projects/hanic/](http://oss.oracle.com/projects/hanic/)

ORACLE

# NIC Bonding cont.

- Solaris
  - IPMP: http://wwws.sun.com/software/solaris/ds/ds-netmultipath/index.html
- HP
  - Auto Port Aggregation  (HPUX): http://www.hp.com/products1/serverconnectivity/adapters/apa_overview.html
  - (Tru64):
- AIX
  - Etherchannel: http://www-1.ibm.com/support/techdocs/atsmastr.nsf/WebIndex/TD101260
- Windows

**ORACLE**

# Shared Storage Configuration

- Configure devices for the Voting Disk and OCR file.
  - Voting Disk >= 20MB, OCR >= 100MB.
  - Use storage mirroring to protect these devices
- Configure shared Storage (for ASM)
  - Use large number of similarly sized "disks"
  - Confirm shared access to storage "disks" from all nodes
  - Use storage mirroring if available
  - Include space for flash recovery area
- Configure IO Multi-pathing
  - ASM must only see a single (virtual) path to the storage
  - Multi-pathing configuration is platform specific (e.g. Powerpath, SecurePath, …)
- Establish file system or location for ORACLE_HOME
  - (And CRS & ASM HOME)
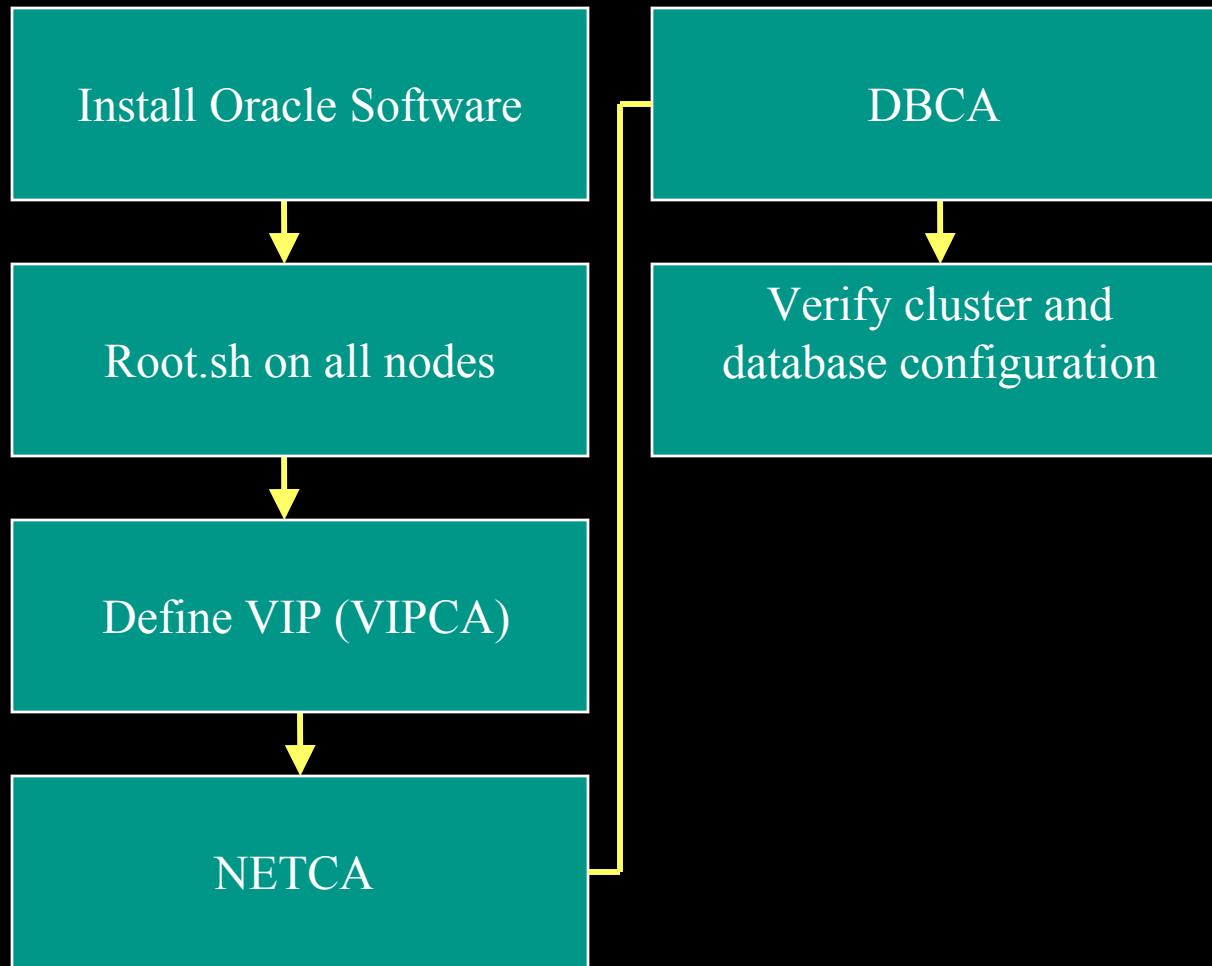
# Installation Flowchart CRS

# Oracle Cluster Manager (CRS) Installation

- CRS is REQUIRED to be installed and running prior to installing 10g RAC.

- CRS must be installed in a different location from the ORACLE_HOME, (e.g. ORA_CRS_HOME).

- Shared Location(s) or devices for the Voting File and OCR file must be available PRIOR to installing CRS.
  - Reinstallation of CRS requires re-initialization of devices, including permissions.

- CRS and RAC require that the private and public network interfaces be configured prior to installing CRS or RAC

- Specify virtual interconnect for CRS communication

# CRS Installation cont.

- Only one set of CRS daemons can be running per RAC node.

- On Unix, the CRS stack is run from entries in /etc/inittab with 'respawn'.

- The supported method to start CRS is booting the machine

- The supported method to stop is shutdown the machine or use "init.crs stop".

ORACLE

# Installation Flowchart Oracle

# Oracle Installation

- The Oracle 10g Installation can be performed after CRS is installed and running on all nodes.

- Start the runInstaller (do not cd in your /mnt/cdrom directory)

- Run root.sh on all nodes
    - Running root.sh on the first node will invoke VIPCA who will configure your Virtual IP's on all nodes
    - After root.sh is finished on the first node start this one after each other on the remaining nodes.

**ORACLE**

# VIP Installation

- The VIP Configuration Assistant (vipca) starts automatically from $ORACLE_HOME/root.sh

- After the welcome screen you have to choose only the public interfaces(s)

- The next screen will ask you for the Virtual IPs for cluster nodes, add your /etc/hosts defined name under IP Alias Name.

  - The VIP must be a DNS known IP address because we use the VIP for the tnsnames connect.

- After finishing this you will see a new VIP interface eg: eth0:1. Use ifconfig (on most platforms) to verify this.

# VIP Installation cont.

- If a cluster is moving to a new datacenter (or subnet) it is necessary to change IPs. The VIP is stored within the OCR and any modification or change to the IP requires additional administrative steps
    - Please see Metalink Note:276434.1 for details

**ORACLE**

# NETCA Best Practices?

- Configure Listeners to listen on the VIP, not on the hostname

- Server side Load balancing configuration recommendations?

- FaN/FCF configuration recommendations?

- Client-side load balancing?

- SQL*Net parameters?? Recv-timeout, send-timeout?

**ORACLE**

# Create RAC database using DBCA

- Set MAXINSTANCES, MAXLOGFILES, MAXLOGMEMBERS, MAXLOGHISTORY, MAXDATAFILES (auto with DBCA)

- Create tablespaces as locally Managed (auto with DBCA)

- Create all tablespaces with ASSM (auto with DBCA)

- Configure automatic UNDO management (auto with DBCA)

- Use SPFILE instead of multiple init.ora's (auto with DBCA)

**ORACLE**

# ASM Disk(group) Best Practices

- ASM configuration performed initially as part of DBCA
- Generally create 2 diskgroups.
  - database area
  - flash recovery area
    - Size dependant on what is stored, and retention period
- Physically separate the database and flashback areas, making sure the two areas do not share the same physical spindles.
- Use diskgroups with large number of similarly sized disks.
- When performing mount operations on diskgroups, it is advisable to mount all required diskgroups at once.
- Make sure disks span several backend disk adapters.
- If mirroring is done in the storage array, set REDUNDANCY=EXTERNAL
- Where possible, use the pseudo devices (multi-path IO)  as the diskstring for ASM.

# ASM File Best Practices

- Use OMF with ASM
- Set db_create_file_dest=+group1
- Create tablespace books;
  - select a.name, f.bytes from v$asm_alias a, v$asm_file f where f.file_number=a.file_number;

```
NAME            BYTES
-----------  ----------
Books.256.1  104857600
```

# ASM File Best Practices

- Use User Templates when necessary.
- User or System templates can be specified in ASM file names for creation
- In ASM instance
  - `Alter diskgroup group1 add template fine attributes (fine unprot);`
- In DB instance
  - `create tablespace tb1 datafile '+group1/tb1(fine)' size 100M;`

# Validate Cluster Configuration

- Query OCR to confirm status of all defined services: crsstat –t
- Use script from Note 259301.1 to improve output formatting/readability

| HA Resource | Target | State |
| --- | --- | --- |
| ora.BCRK.BCRK1.inst | ONLINE | ONLINE on sunblade-25 |
| ora.BCRK.BCRK2.inst | ONLINE | ONLINE on sunblade-26 |
| ora.BCRK.db | ONLINE | ONLINE on sunblade-25 |
| ora.sunblade-25.ASM1.asm | ONLINE | ONLINE on sunblade-25 |
| ora.sunblade-25.LISTENER_SUNBLADE-25.lsnr | ONLINE | ONLINE on sunblade-25 |
| ora.sunblade-25.gsd | ONLINE | ONLINE on sunblade-25 |
| ora.sunblade-25.ons | ONLINE | ONLINE on sunblade-25 |
| ora.sunblade-25.vip | ONLINE | ONLINE on sunblade-25 |
| ora.sunblade-26.ASM2.asm | ONLINE | ONLINE on sunblade-26 |
| ora.sunblade-26.LISTENER_SUNBLADE-26.lsnr | ONLINE | ONLINE on sunblade-26 |
| ora.sunblade-26.gsd | ONLINE | ONLINE on sunblade-26 |
| ora.sunblade-26.ons | ONLINE | ONLINE on sunblade-26 |
| ora.sunblade-26.vip | ONLINE | ONLINE on sunblade-26 |

# Validate RAC Configuration

- Instances running on all nodes
    SQL> select * from gv$instance

- RAC communicating over the private Interconnect
    SQL> oradebug setmypid
    SQL> oradebug ipc
    SQL> oradebug tracefile_name
    /home/oracle/admin/RAC92_1/udump/rac92_1_ora_1343841.trc
    - Check trace file in the user_dump_dest:
        SSKGXPT 0x2ab25bc flags        info for network 0
            socket no 10   IP **10.0.0.1**   UDP 49197
            sflags SSKGXPT_UP
            info for network 1
            socket no 0    IP 0.0.0.0      UDP 0
            sflags SSKGXPT_DOWN

# Validate RAC Configuration

- RAC is using desired IPC protocol: Check Alert.log

  cluster interconnect IPC version:Oracle UDP/IP

  IPC Vendor 1 proto 2 Version 1.0

  PMON started with pid=2

- Use cluster_interconnects only if necessary

  – RAC will use the same "virtual" interconnect selected during CRS install

  – To check which interconnect and is used and where it came from use "select * from x$ksxpia;"

```
ADDR                    INDX      INST_ID P PICK NAME_KSXPIA      IP_KSXPIA
--------------- ---------- ---------- - ---- --------------- ---------
00000003936B8580         0          1   OCR  eth1             10.0.0.1

Pick: OCR … Oracle Clusterware
      OSD … Operating System dependent
       CI … indicates that the init.ora parameter
            cluster_interconnects is specified
```

# Post Installation

- Enable asynchronous I/O if available
  - cd $ORACLE_HOME/rdbms/lib; make -f ins_rdbms.mk async_on ioracle
- Adjust UDP send / receive buffer size to 256K (Linux only)
- If Buffer Cache > 1.7GB required, use 64-bit platform.

**ORACLE**

# Optimize Instance Recovery

- Set fast_start_mttr_target
  - 60 < fsmttr < 300 is a good starting point
  - Balance of performance vs. availability
- Size the buffer cache for single pass recovery.
- Ensure asynch I/O is used.
- **Follow configuration best practices as documented in Oracle® High Availability Architecture and Best Practices 10*g* Release 1 (10.1)**

# SRVCTL

- SRVCTL is a very powerful tool
- SRVCTL uses information from the OCR file
- GSD in 10g is running just for compatibility to serve 9i clients if 9i and 10g is running on the same cluster.
- srvctl status nodeapps -n <nodename> will show all services running on a node
  - SRVCTL commands are documented in Appendix B of the RAC Admin Guide at:
    http://download-west.oracle.com/docs/cd/B13789_01/rac.101/b10765/toc.htm

ORACLE

# Application Considerations
# FCF vs. TAF

- Connection Retries:
  - FCF allows retry at the Application level, TAF retries occur at the OCI/Net layer. Application layer (Example: EJB Container) fully controls retries

- Integrated with the Connection Cache:
  - FCF works in conjunction with the Implicit Connection Cache, and has complete control over connections managed by the cache

- RAC Events Based:
  - FCF is a RAC event based mechanism. This is more efficient than detecting failures of network calls

- Load Balancing Support:
  - FCF supports UP event Load Balancing of connections across active RAC instances – start and UP
  - Work requests are distributed across RAC

**ORACLE**

# Applications Waste Time

| Connect | SQL issue | Blocked in R/W | Processing last result |
|---|---|---|---|
| active | active | wait | wait |
| tcp_ip_cinterval | tcp_ip_interval | tcp_ip_keepalive | - |
| VIP | VIP | out of band event - FAN | out of band event - FAN |

# What is FaN?

- Fast Application Notification  (FaN) is RAC HA notification mechanism which let applications know about service & node events (UP or DOWN events)

- Fast Connection Failover (FCF) is a mechanism of 10g JDBC which uses FAN

- Enable it, and Forget it. Works transparently by receiving asynchronous events from the RAC database

**ORACLE**

# How does Fast Connection Failover use FAN?

- FCF is a subscriber of FAN, where
  - instance UP event - leverages FaN to load balance connections across the existing and new instances
  - node/instance DOWN event - cleans up the connection cache (remove invalid connections)
- iAS 10.1.3 will integrate JDBC 10g
- Query/Operations retries are up to the application/containers not FCF.
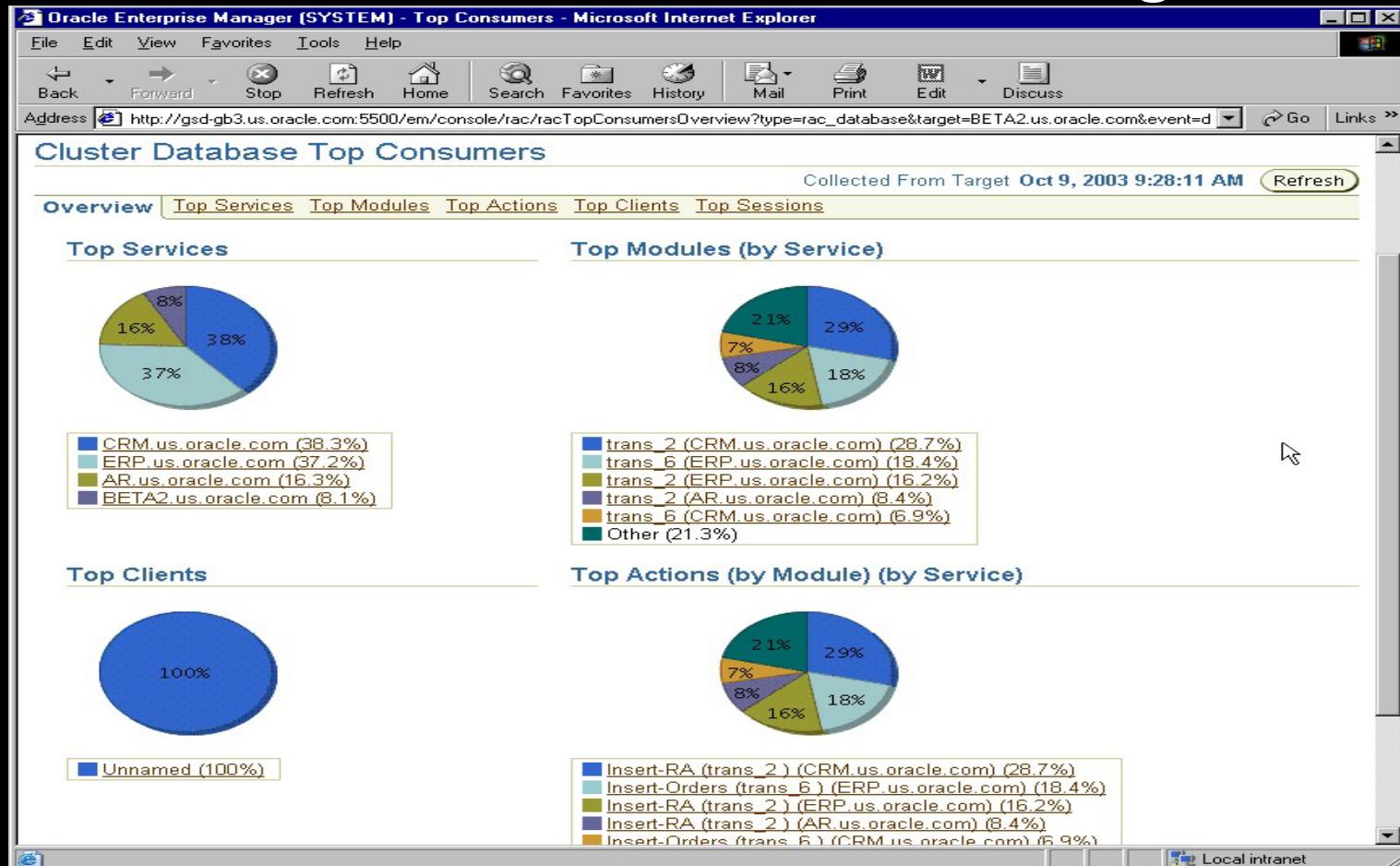
ORACLE

# What is a Service?

- In Oracle10g services are built into the database.
- Divides work into logical workloads which share common functions, service level thresholds, priority & resource needs.
- Examples:
    - OLTP & Batch
    - ERP, CRM, HR, Email
    - DW & OLTP
    - Affinity Group 1,2,3,4,5,6,7,8,9,10

**ORACLE**

# Take Advantage of 10g Services

- Easy to setup, configure then connect by service
- Benefits
  - Availability
    - Services has a defined Topology & automatic recovery
    - Callouts as services come up and down
  - Performance
    - A new level for performance tuning
    - Workload are routed transparently
    - Alerts & actions when performance goals are violated
    - Natural support for mixed workloads and mixed size nodes
  - Manageability
    - Each workload is managed in isolation
    - Prioritization & Resource Management

# Services in Enterprise Manager
## *Critical Tool for Performance Tuning*



*More details on Services provided in a separate Web Seminar*

ORACLE

# Application Considerations Configuration

- Plan your services
  - application to service, data range to service
  - global name, HA configuration, priority, response time
- Use service: not SID, not Instance, not Host
  - Use service to connect
  - Use virtual IP for database access
  - Use cluster alias to eliminate address lists.
- Use service for jobs and PQ.

# Application Considerations Runtime

- Make applications measurable
  - instrument with MODULE and ACTION
  - use the DBMS_MONITOR to gather statistics
- For priorities – use resource manger
- For load balancing
  - use CLB to balance connections by service.
  - use service metrics to "deal requests" from mid-tier connection pools by service.

# Application Considerations Recovery

- Use JDBC connection pools for fast failover.
    - Surviving sessions continue *FAST*.
    - Interrupted sessions detect the error *FAST*.
- Use TAF callbacks to trap and handle errors.
- Use HA callouts/events (up, down, not restarting) to notify the application to take appropriate action.
    - Save and recall non-transactional state.
    - Check transaction outcome and resubmit.

# Application Deployment

- Same guidelines as single instance
  - SQL Tuning
  - Sequence Caching
  - Partition large objects
  - Use different block sizes
  - Tune instance recovery
  - Avoid DDL
  - Use LMT's and ASSM

# Agenda

- Planning Best Practices
  - Architecture
  - Expectation setting
  - Objectives and success criteria
  - Project plan

- Implementation Best Practices
  - Infrastructure considerations
  - Installation/configuration
  - Database creation
  - Application considerations

- Operational Best Practices
  - Backup & Recovery
  - Performance Monitoring and Tuning
  - Production Migration

**ORACLE**

# Operations

- Same DBA procedures as single instance, with some minor, mostly mechanical differences.
- Managing the Oracle environment
  - Starting/stopping Oracle cluster stack with boot/reboot server
  - Managing multiple redo log threads
- Startup and shutdown of the database
  - Use Grid Control
- Backup and recovery
- Performance Monitoring and Tuning
- Production migration

**ORACLE**

# Operations: Backup & Recovery

- RMAN is the most efficient option for Backup & Recovery
  - Managing the snapshot control file location.
  - Managing the control file autobackup feature.
  - Managing archived logs in RAC – choose proper archiving scheme.
  - Node Affinity Awareness

- RMAN and Oracle Net in RAC apply
  - you cannot specify a net service name that uses Oracle Net features to distribute RMAN connections to more than one instance.

- Oracle Enterprise Manager
  - GUI interface to Recovery Manager

# Backup & Recovery

- Use RMAN
  - Only option to backup and restore ASM files
- Use Grid Control
  - GUI interface to RMAN
- Use 10g Flash Recovery Area for backups and archive logs
  - On ASM and available to all instances

**ORACLE**

# Performance Monitoring and Tuning

- Tune first for single instance 10g
- Use ADDM and AWR
- Oracle Performance Manager
- RAC-specific views
- Supplement with scripts/tracing
  - Monitor V$SESSION_WAIT to see which blocks are involved in wait events
  - Trace events like 10046/8 can provide additional wait event details
  - Monitor Alert logs and trace files, as on single instance
- Supplement with System-level monitoring
  - CPU utilization never 100%
  - I/O service times never > acceptable thresholds
  - CPU run queues at optimal levels
- Note that in 10g, performance statistics are message/time based, as opposed to event-based in 9i

ORACLE

# Performance Monitoring and Tuning

- Obvious application deficiency on a single node can't be solved by multiple nodes.
    - Single points of contention.
    - Not scalable on SMP
    - I/O bound on single instance DB
- Tuning on single instance DB to ensure applications scalable first
    - Identify/tune contention using v$segment_statistics to identify objects involved
    - Concentrate on the top wait events if majority of time is spent waiting
    - Concentrate on bad SQL if CPU bound
- Maintain a balanced load on underlying systems (DB, OS, storage subsystem, etc. )
    - Excessive load on individual components can invoke aberrant behaviour.

**ORACLE**

# Performance Monitoring / Tuning

- Deciding if RAC is the performance bottleneck
  - "Cluster" wait event class
  - Amount of Cross Instance Traffic
    - Type of requests
    - Type of blocks
  - Latency
    - Block receive time
    - buffer size factor
    - bandwidth factor

# Avoid false node evictions

- May get 'heart beat' failures if critical processes are unable to respond quickly
  - Enable real time priority for LMS
  - Do not run system at 100% CPU over long period
  - Ensure good I/O response times for control file and voting disk

# Production Migration

- Adhere to strong Systems Life Cycle Disciplines
  - Comprehensive test plans (functional and stress)
  - Rehearsed production migration plan
  - **Change Control**
    - **Separate environments for Dev, Test, QA/UAT, Production**
    - **System AND application change control**
    - **Log changes to spfile**
  - Backup and recovery procedures
  - Patchset maintenance
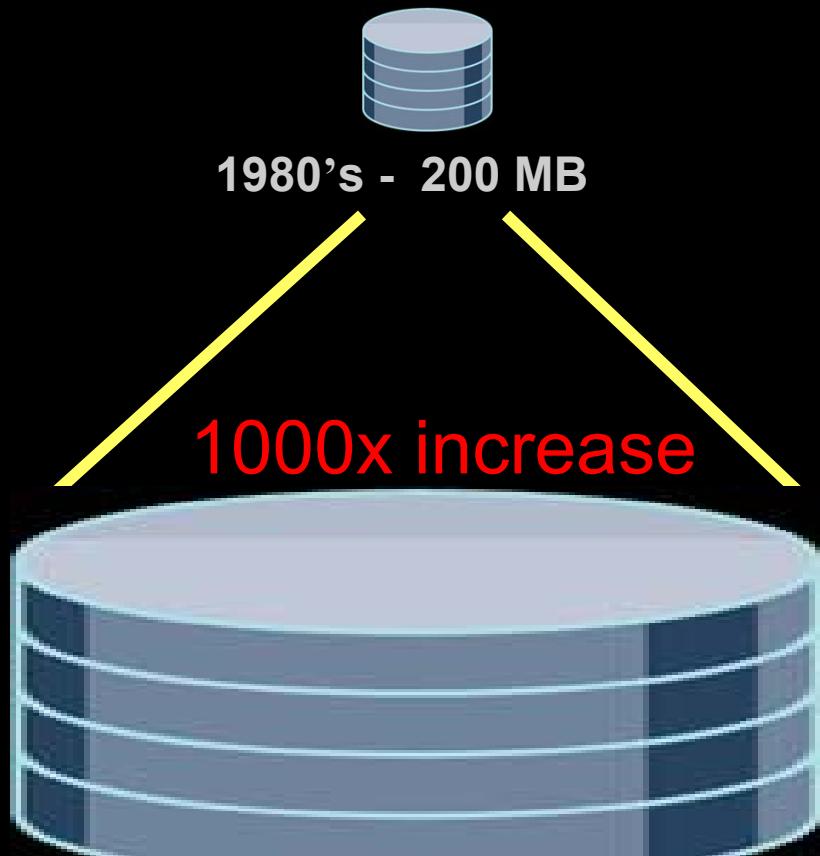  - Security controls
  - Support Procedures

# New World:
# Disk Based Data Recovery
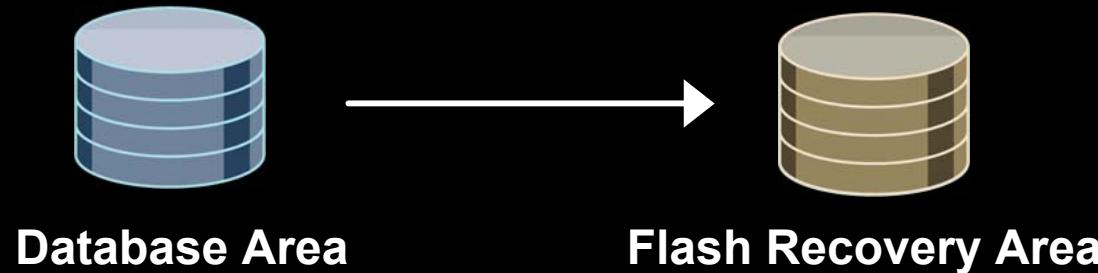
**1980's - 200 MB**
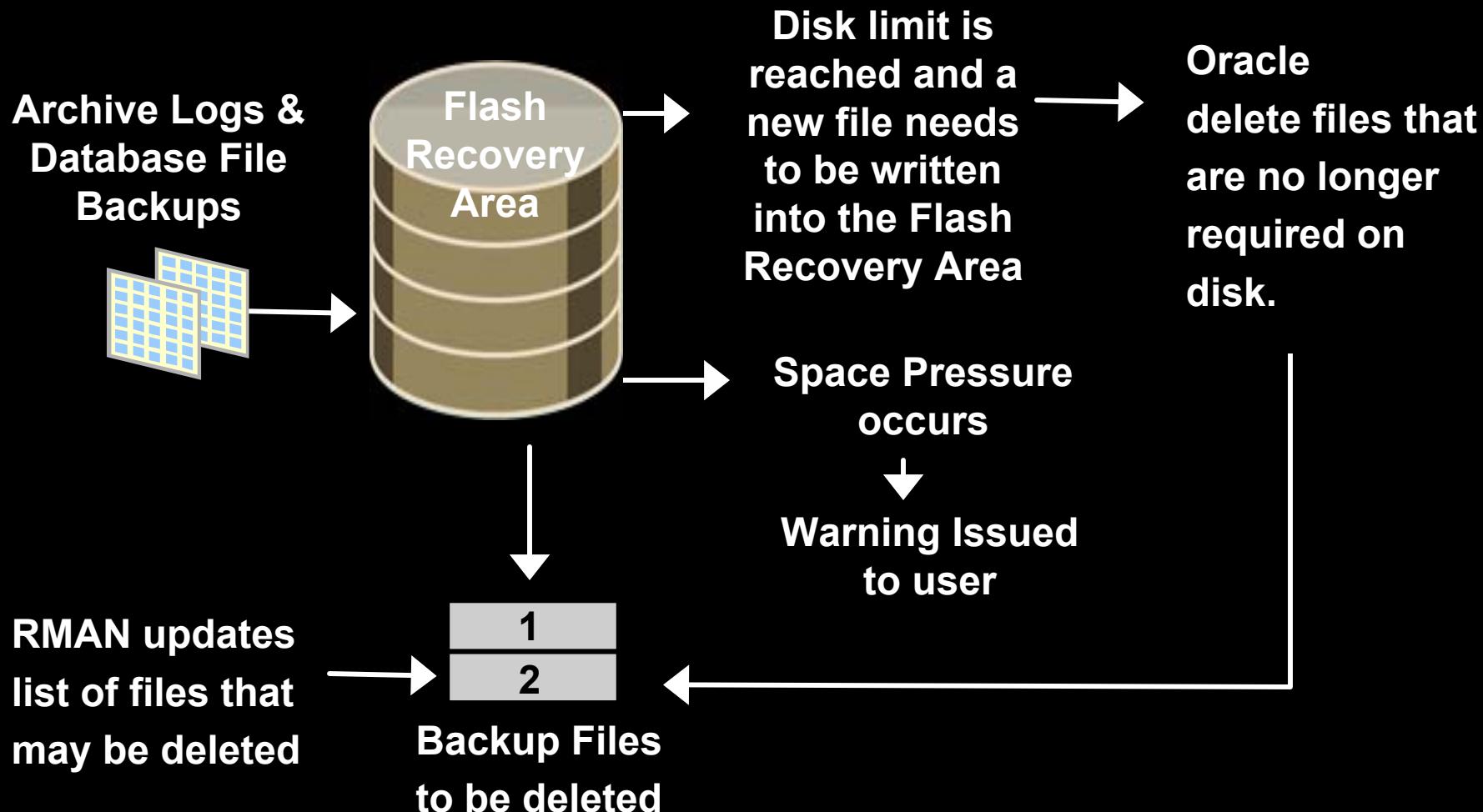
## 1000x increase

**2000's - 200 GB**

- **Disk economics are close to tape**
- **Disk is better than tape**
  - Random access to any data
- **We rearchitected our recovery strategy to take advantage of these economics**
  - Random access allows us to backup and recover just the <u>changes</u> to the database
- **Backup and Recovery goes from hours to minutes**

# Flash Recovery Area

- **Unified storage location for all recovery files and recovery related activities in an Oracle Database.**
    - **Centralized location for control files, online redo logs, archive logs, flashback logs, backups**
    - **A flash recovery area can be defined as a directory, file system, or ASM disk group**
    - **A single recovery area can be shared by more than one database**
- **Minimize the number of initialization parameters to set when you create a database**
    - **Define a database area and flash recovery area location**
    - **Oracle creates and manages all files using OMF**



**Database Area**        **Flash Recovery Area**

# Flash Recovery Area Space Management

**Archive Logs & Database File Backups**

**Flash Recovery Area**

**Disk limit is reached and a new file needs to be written into the Flash Recovery Area**

**Oracle delete files that are no longer required on disk.**

**Space Pressure occurs**

**Warning Issued to user**

**RMAN updates list of files that may be deleted**

| 1 |
|---|
| 2 |

**Backup Files to be deleted**

# Benefits to Using a Flash Recovery Area

- <u>Unifies</u> the storage location of related recovery files

- <u>Manages</u> the disk space allocated for recovery files automatically

- <u>Simplifies</u> database administrator tasks

- Much <u>faster backup</u>

- Much <u>faster restore</u>

- Much <u>more reliable</u> due to inherent reliability of disks

**ORACLE**