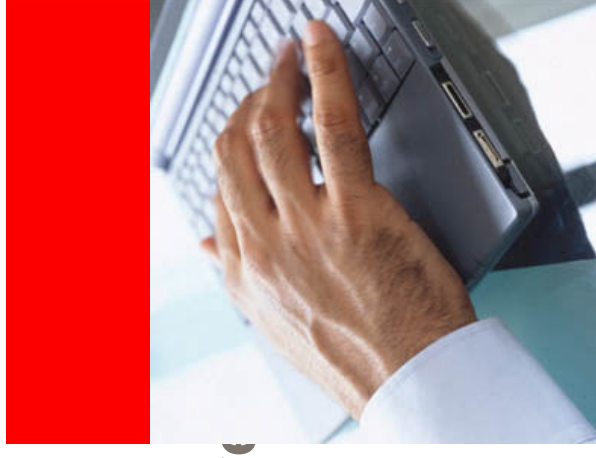# ORACLE®

## RAC Performance Tuning Best Practices

The following is intended to outline our general product direction. It is intended for information purposes only, and may not be incorporated into any contract. It is not a commitment to deliver any material, code, or functionality, and should not be relied upon in making purchasing decisions. The development, release, and timing of any features or functionality described for Oracle's products remains at the sole discretion of Oracle.

ORACLE®

# Agenda

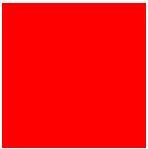## Practical RAC Performance Analysis Revi

- RAC Architecture Overview
- Common Problems and Symptoms
- Application and Database Design
- Diagnostics and Problem Determination
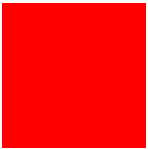- Summary: Practical Performance Analysis
- Appendix

# OBJECTIVE

- Realize that RAC performance does not requires "Black Magic"

- General system and SQL analysis and tuning experience is practically sufficient for RAC

- Problems can be identified with a minimum of metrics and effort

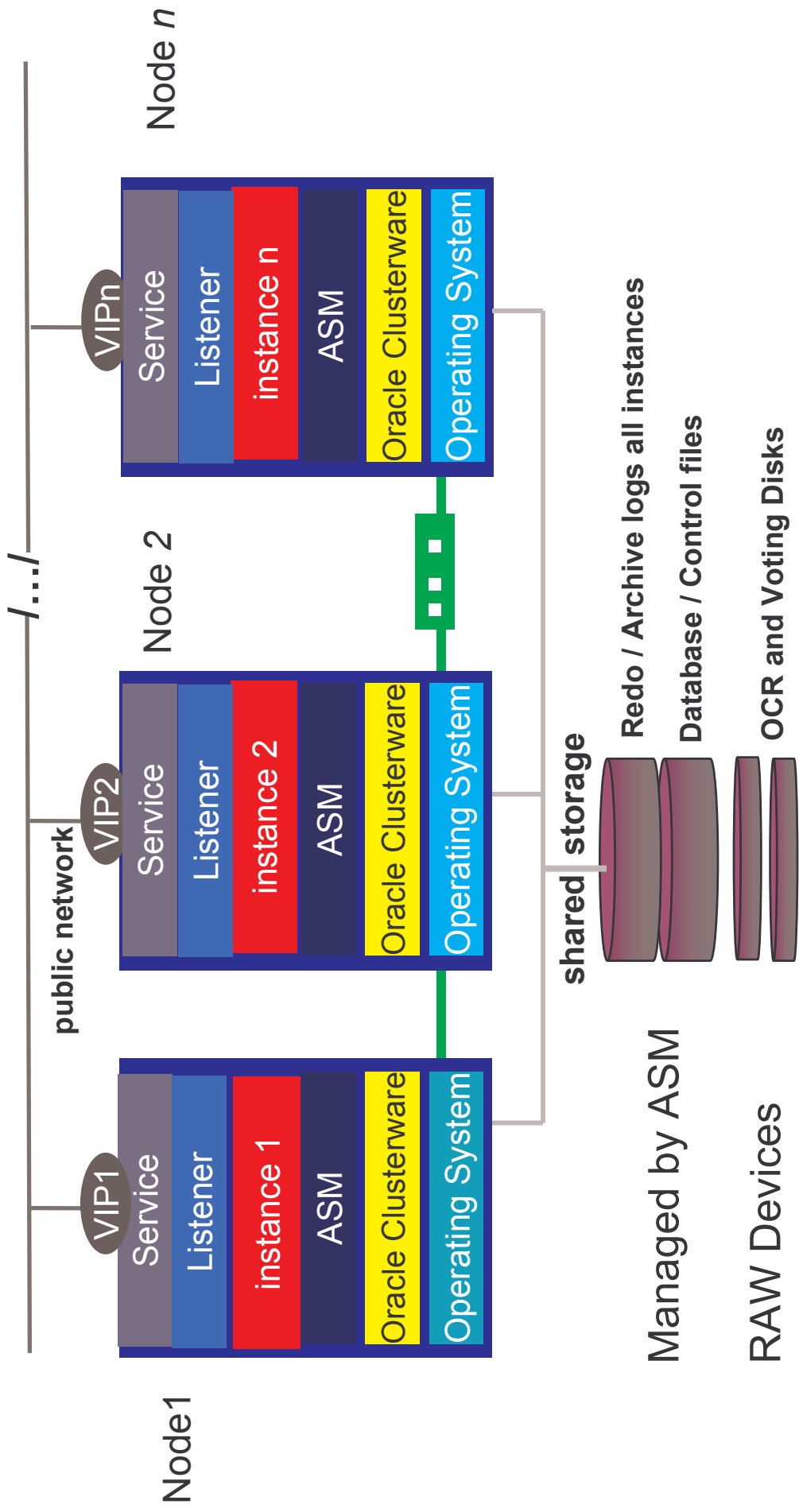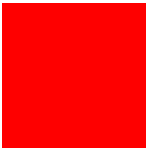- Diagnostics framework and Advisories are efficient
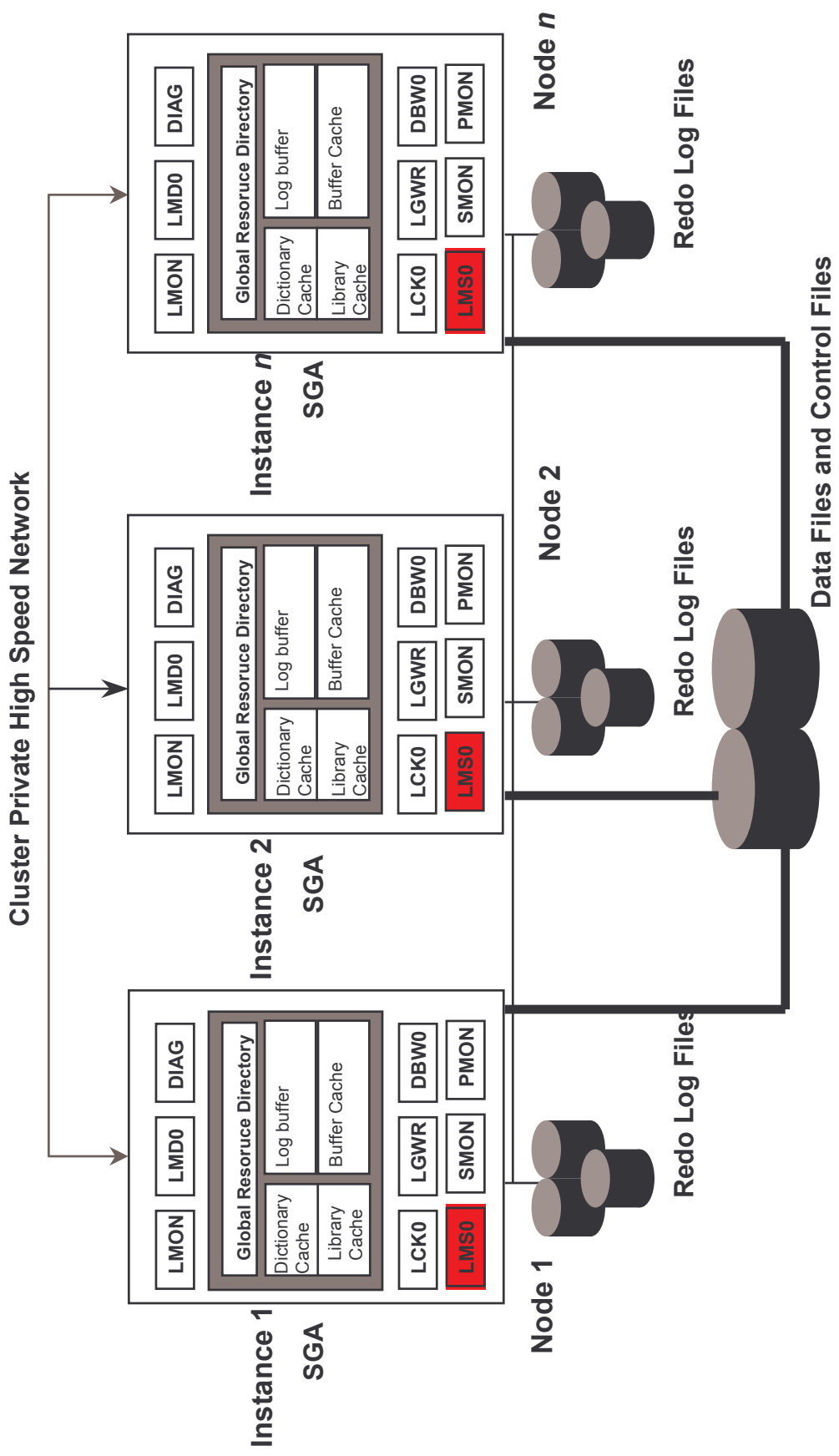
ORACLE

# RAC Architecture Overview

# RAC 10g Architecture

# Under the Covers



**Cluster Private High Speed Network**

**Instance 1 SGA**

LMON | LMD0 | DIAG

Global Resoruce Directory
Dictionary Cache | Log buffer
Library Cache | Buffer Cache

LCK0 | LGWR | DBW0
LMS0 | SMON | PMON

**Node 1**

Redo Log Files

**Instance 2 SGA**

LMON | LMD0 | DIAG

Global Resoruce Directory
Dictionary Cache | Log buffer
Library Cache | Buffer Cache

LCK0 | LGWR | DBW0
LMS0 | SMON | PMON

**Node 2**

Redo Log Files

**Instance _n_ SGA**

LMON | LMD0 | DIAG

Global Resoruce Directory
Dictionary Cache | Log buffer
Library Cache | Buffer Cache

LCK0 | LGWR | DBW0
LMS0 | SMON | PMON

**Node _n_**

Redo Log Files

**Data Files and Control Files**

ORACLE
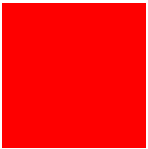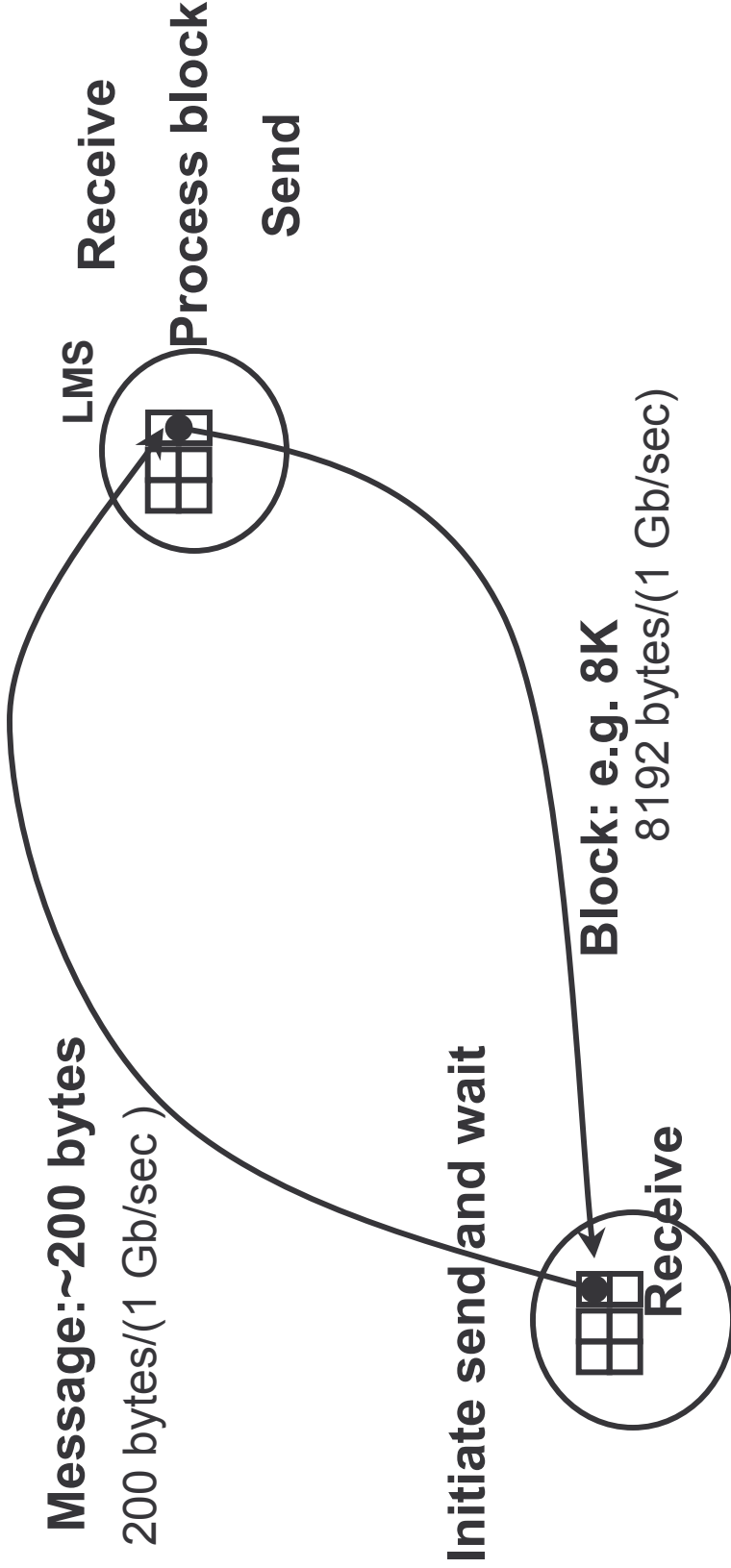
# Global Cache Service (GCS)

- Guarantees cache coherency

- Manages caching of shared data via Cache Fusion

- Minimizes disk access to data which is not in local cache by remotely transferring blocks

- Implements fast direct memory access over high-speed interconnects for all data blocks and types

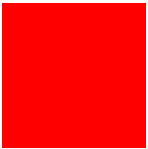- Uses an efficient and scalable messaging protocol

ORACLE

# GCS Processing

**Message:~200 bytes**
200 bytes/(1 Gb/sec )

**LMS**  **Receive**

**Process block**

**Send**

**Block: e.g. 8K**
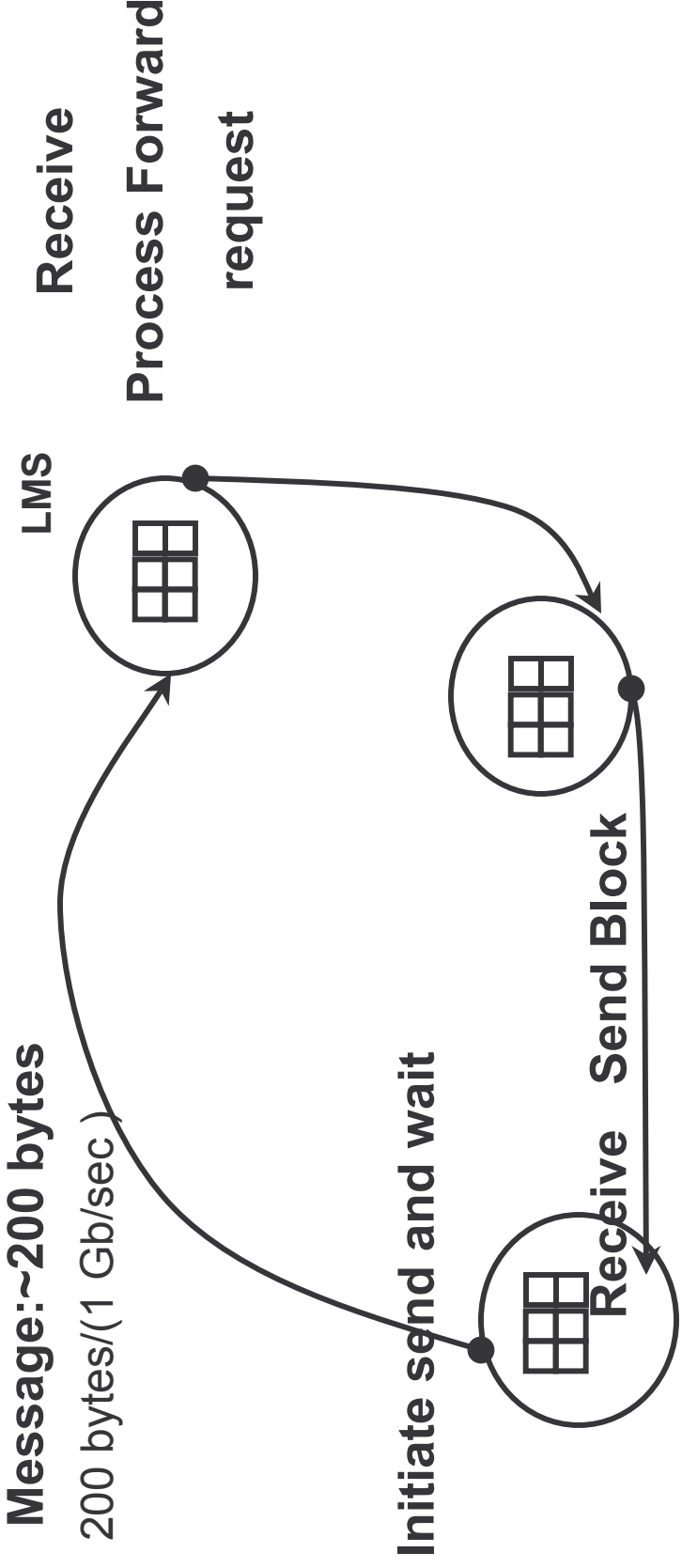8192 bytes/(1 Gb/sec )

**Initiate send and wait**

**Receive**

Total access time: e.g. ~360 microseconds (UDP over GBE)
Network propagation delay ( "wire time" ) is a minor factor for roundtrip time
( approx.: 6% , vs. 52% in OS and network stack )

# GCS Processing

**Message:~200 bytes**

200 bytes/(1 Gb/sec )

**Receive**

**Process Forward**

**request**

LMS

Initiate send and wait

**Receive   Send Block**

Network propagation delay ( "wire time" )  is a minor factor for roundtrip time

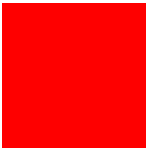( approx.: 6% , vs. 52% in OS and network stack )

ORACLE®

# Block Transfer Time

Determined by

- Network Transmit Time – aka wire speed
- HOST CPU
  - Send/Receive Network packet processing
  - LMS processing
- Operating System scheduling
- LMS load
- Interconnect stability

- Oracle statistics report Round-trip Time

# Block Transfer Latency

- ~300 microseconds is lowest measured with UDP over Gigabit Ethernet and 2K blocks

- ~ 120 microseconds is lowest measured with RDS over Infiniband and 2K blocks

| Block size RT (ms) | 2K | 4K | 8K | 16K |
|---|---|---|---|---|
| UDP/GE | 0.30 | 0.31 | 0.36 | 0.46 |
| RDS/IB | 0.12 | 0.13 | 0.16 | 0.20 |

ORACLE

# Infrastructure: Private Interconnect

- Network between the nodes of a RAC cluster MUST be private

- Supported links: GbE,  IB ( IPoIB: 10.2 )

- Supported transport protocols: UDP, RDS (10.3)

- Use multiple or dual-ported NICs for redundancy and increase bandwidth with NIC bonding

- Large ( Jumbo ) Frames for GbE recommended

# Infrastructure: Interconnect Bandwidth

- Bandwidth requirements depend on
  - CPU power per cluster node
  - Application-driven data access frequency
  - Number of nodes and size of the working set
  - Data distribution between PQ slaves
- Typical utilization approx. 10-30% in OLTP
  - 10000-12000 8K blocks per sec to saturate 1 x Gb Ethernet ( 75-80% of theoretical bandwidth )
- Multiple NICs generally not required for performance and scalability

ORACLE

# Infrastructure: IPC configuration

- Settings:
  - Socket receive buffers ( 256 KB – 1MB )
  - Negotiated top bit rate and full duplex mode
  - NIC ring buffers
  - Ethernet flow control settings
- Verify your setup:
  - CVU does checking
  - Load testing eliminates potential for problems

# Infrastructure: Operating System

- Remote Block access latencies increase when CPU(s) busy and run queues are long

- Immediate LMS scheduling is critical for predictable block access latencies when CPU > 80% busy

- Real Time or fixed priority for LMS is supported

  - Implemented by default with 10.2

ORACLE

# Infrastructure: IO capacity

- Disk storage is shared by all nodes, i.e the aggregate IO rate is important
- Log file IO latency can be important for block transfers
- Parallel Execution across cluster nodes requires a well-scalable IO subsystem
  - Disk configuration needs to be responsive and scalable
  - Get I/O baseline with ORION

# Common Problems and Symptoms

# Misconfigured or Faulty Interconnect Can Cause:

- Dropped packets/fragments
- Buffer overflows
- Packet reassembly failures or timeouts
- Ethernet Flow control issues
- TX/RX errors

"gc lost blocks" responsible for large no of escalations

ORACLE

# "Lost Blocks": NIC Receive Errors

**Db_block_size = 8K**

```
ifconfig -a:

eth0 Link encap:Ethernet  HWaddr 00:0B:DB:4B:A2:04

     inet addr:130.35.25.110  Bcast:130.35.27.255  Mask:255.255.252.0

     UP BROADCAST RUNNING MULTICAST  MTU:1500  Metric:1

     RX packets:21721236 errors:135 dropped:0 overruns:0 frame:95

     TX packets:273120 errors:0 dropped:0 overruns:0 carrier:0

     ...
```

# "Lost Blocks": IP Packet Reassembly Failures

netstat –s

Ip:
848847742 total packets received

...
**1201 fragments dropped after timeout**

...
**3384 packet reassembles failed**

ORACLE

# Finding a Problem with the Interconnect or IPC

```
Top 5 Timed Events
~~~~~~~~~~~~~~~~~~~
                                            Avg    %Total
                                            wait   Call
Event              Waits     Time(s)  (ms)  Time   Wait Class
-------------------------------------------------------------
log file sync      286,038   49,872   174   41.7   Commit
gc buffer busy     177,315   29,021   164   24.3   Cluster
gc cr block busy   110,348    5,703    52    4.8   Cluster
gc cr block lost     4,272    4,953  1159    4.1   Cluster
cr request retry     6,316    4,668   739    3.9   Other
```

*Should never be here*

# Impact of IO capacity issues or bad SQL execution on RAC

- Log flush IO delays can cause "busy" buffers
- "Bad" queries on one node can saturate the link
- I/O is issued from ALL nodes to shared storage ( beware of one-node "myopia" )

Cluster-wide impact of IO or query plan issues responsible substantial no of escalations

ORACLE®

# Cluster-Wide IO Impact

## Node 1

```
Top 5 Timed Events
~~~~~~~~~~~~~~~~~~~
                                           Avg  %Total
                                          wait  Call
Event                  Waits   Time(s)  (ms)  Time
-------------------   -------  -------  ----  ------
log file sync         286,038  49,872   174    41.7
gc buffer busy        177,315  29,021   164    24.3
gc cr block busy      110,348   5,703    52     4.8
```

## Node 2

```
Load Profile                    Per Second
~~~~~~~~~~~                      ----------
Redo size:                       40,982.21
Logical reads:                   81,652.41
Physical reads:                  51,193.37
```

ORACLE

# IO and bad SQL problem fixed

```
Top 5 Timed Events
~~~~~~~~~~~~~~~~~~
                                                 Avg   %Total
                                                 wait   Call
Event                     Waits      Time (s)    (ms)   Time   Wait Class
-------------             -------    --------    ----   -----  ------------
CPU time                  4,580         65.4
log file sync             276,281       1,501      5    21.4   Commit
log file parallel write   298,045         923      3    13.2   System I/O
gc current block 3-way    605,628         631      1     9.0   Cluster
gc cr block 3-way         514,218         533      1     7.6   Cluster
```

ORACLE

# CPU Saturation or Memory Depletion

```
Top 5 Timed Events
~~~~~~~~~~~~~~~~~~
                                  Avg  %Total
                            wait  Call
Event           Waits    Time(s) (ms) Time  Wait Class
--------------- -------  ------- ----- ----- -----------
db file sequential 1,312,840 21,590 16  21.8  User I/O
read

gc current block   275,004 21,054  77  21.3  Cluster
congested

gc cr grant congested 177,044 13,495  76  13.6  Cluster

gc current block  1,192,113  9,931   8  10.0  Cluster
2-way

gc cr block congested  85,975  8,917 104   9.0  Cluster
```

*"Congested": LMS could not process block transfer request fast enough*
*Cause      : Long run-queues and paging on the cluster nodes*

ORACLE

# Health Check

Look for:

- High impact of "lost blocks" , e.g.

  gc cr block lost          1159

- IO capacity saturation , e.g.

  gc cr block busy          52 ms

- Overload and memory depletion, e.g

  gc current block congested          14 ms

ORACLE

# Application and Database Design

# General Principles

- No fundamentally different design and coding practices for RAC
- Badly tuned SQL and schema will not run better
- Serializing contention makes applications less scalable
- Standard SQL solves > 80% of performance problems
- Follow RAC Best Practices – Accumulation of Real-world knowledge

# Scalability Pitfalls

- Serializing contention on a small set of data/index blocks
  - monotonically increasing key
  - frequent updates of small cached tables
  - segment without ASSM or Free List Group (FLG)
- Full table scans
- Frequent hard parsing
- Concurrent DDL ( e.g. truncate/drop )

# Index Block Contention: Optimal Design

- Monotonically increasing sequence numbers
  - Randomize or cache
  - Large ORACLE sequence number caches
- Hash or range partitioning
  - Local indexes

# Data Block Contention: Optimal Design

- Small tables with high row density and frequent updates and reads can become "globally hot" with serialization e.g.
  - Queue tables
  - session/job status tables
  - last trade lookup tables
- Higher PCTFREE for table reduces # of rows per block

# Large Contiguous Scans

- Query Tuning
- Use parallel execution
  - Intra- or inter instance parallelism
  - Direct reads
  - GCS messaging minimal

# Health Check

Look for:

- Indexes with right-growing characteristics
  - Eliminate indexes which are not needed
- Frequent updated and reads of "small" tables
  - "small"=fits into a single buffer cache
- SQL which scans large amount of data
  - Bad execution plan
  - More efficient when parallelized

# Diagnostics and Problem Determination

# Performance Checks and Diagnosis

- Traditionally done via AWR or Statspack reports
- "Time-based" paradigm, i.e. identify which events consume the highest proportion of the database time
- Global cache ( "gc" ) events are typical for RAC
- Drill-down to SQL and Segment Statistics

ORACLE

# Event Statistics to Drive Analysis

- Global cache ("gc") events and statistics
  - Indicate that Oracle searches the cache hierarchy to find data fast
  - as "normal" as an IO ( e.g. db file sequential read )
- GC events tagged as "busy" or "congested" consuming a significant amount of database time should be investigated
  - At first, assume a load or IO problem on one or several of the cluster nodes

ORACLE

# Global Cache Event Semantics

All Global Cache Events will follow the following format:

GC ...

- CR, current
  - Buffer requests and received for read or write
- block, grant
  - Received block or grant to read from disk
- 2-way, 3-way
  - Immediate response to remote request after N-hops
- busy
  - Block or grant was held up because of contention
- congested
  - Block or grant was delayed because LMS was busy or could  not get the CPU

# "Normal" Global Cache Access Statistics

```
Top 5 Timed Events
~~~~~~~~~~~~~~~~~~
                                               Avg    %Total
                                               wait   Call
Event               Waits        Time(s)       (ms)   Time    Wait Class
------              -----        -------       ----   ----    ----------
CPU time                          4,580        65.4

log file sync       276,281       1,501        5      21.4    Commit

log file parallel   298,045       923          3      13.2    System I/O
write

gc current block    605,628       631          1      9.0     Cluster
3-way

gc cr block 3-way   514,218       533          1      7.6     Cluster
```

*Avg latency is 1 ms or less*

*Reads from remote cache instead of disk*

ORACLE

# "Abnormal" Global Cache Statistics

```
Top 5 Timed Events                    Avg  %Total
~~~~~~~~~~~~~~~~~~~                    wait  Call
Event           Waits    Time(s)  (ms)  Time  Wait Class
--------------- -------- -------- ----- ----- -----------
log file sync   286,038  49,872   174   41.7  Commit
gc buffer busy  177,315  29,021   164   24.3  Cluster
gc cr block busy 110,348  5,703    52    4.8   Cluster
```

*"busy" indicates contention*

*Avg time is too high*

# Checklist for the Performance Analyst ( AWR based )

- Check where most of the time in the database is spend ("Top 5" )
- Check whether gc events are "busy", "congested"
- Check the avg wait time
- Drill down
  - SQL with highest cluster wait time
  - Segment Statistics with highest block transfers

ORACLE®

# Drill-down: An IO capacity problem

```
Top 5 Timed Events

                                                  Avg  %Total
                                                  wait Call
Event                     Waits       Time(s)     (ms) Time  Wait Class
-----                     -----       -------     ---- ----- ----------
db file scattered read    3,747,683   368,301     98   33.3  User I/O
gc buffer busy            3,376,228   233,632     69   21.1  Cluster
db file parallel read     1,552,284   225,218     145  20.4  User I/O
gc cr multi block         35,588,800  101,888     3    9.2   Cluster
request
read by other session     1,263,599   82,915      66   7.5   User I/O
```

*IO contention*

*Symptom of Full Table Scans*

# Drill-down: SQL Statements

*"Culprit": Query that overwhelms IO subsystem on one node*

| Physical Reads | Executions | per Exec | %Total |
|---|---|---|---|
| ------------- | ------------- | ------------- | ------- |
| 182,977,469 | 1,055 | 173,438.4 | 99.3 |

SELECT SHELL FROM ES_SHELL WHERE MSG_ID = :msg_id ORDER BY ORDER_NO ASC

*The same query reads from the interconnect:*

| Cluster Wait Time (s) | CWT % of Elapsd Tim | CPU Time(s) | Executions |
|---|---|---|---|
| ------------- | ------------- | ------------- | ------------- |
| 341,080.54 | 31.2 | 17,495.38 | 1,055 |

SELECT SHELL FROM ES_SHELL WHERE MSG_ID = :msg_id ORDER BY ORDER_NO ASC

# Drill-Down: Top Segments

| Tablespace Name | Object Name | Subobject Name | Obj. Type | GC Buffer Busy | % of Capture |
|---|---|---|---|---|---|
| ESSMLTBL | ES_SHELL | SYS_P537 | TABLE | 311,966 | 9.91 |
| ESSMLTBL | ES_SHELL | SYS_P538 | TABLE | 277,035 | 8.80 |
| ESSMLTBL | ES_SHELL | SYS_P527 | TABLE | 239,294 | 7.60 |
| ... | | | | | |

Apart from being the table with the highest IO demand it was the table with the highest number of block transfers AND global serialization

# … and now for something different:
# Automated Performance Analysis

# Impact of RAC Findings

# Automated Findings and Actions: Interconnect



**Finding**

**Action**

# Automated Findings and Actions: Block Contention

# Automated Findings and Actions: SQL

# Automated SQL Drill-Down

**SQL Text**

**Per SQL Statistics Over Time**

# Summary: Practical Performance Analysis

# Diagnostics Flow

- Start with simple validations :
  - Private Interconnect used ?
  - Lost blocks and failures ?
  - Load and load distribution issues ?
- Check avg latencies, busy, congested events and their significance
- Check OS statistics ( CPU, disk , virtual memory )
- Identify SQL and Segments

**MOST OF THE TIME, A PERFORMANCE PROBLEM IS NOT A RAC PROBLEM**

ORACLE

# Actions

- Interconnect issues must be fixed first
- If IO wait time is dominant , fix IO issues
  - At this point, performance may already be good
- Fix "bad" plans
- Fix serialization
- Fix schema

ORACLE

# Checklist for Practical Performance Analysis

- ADDM provides RAC performance analysis of significant metrics and statistics
  - ADDM findings should always be studied first
  - It provides detailed findings for SQL, segments and blocks
- AWR for detailed statistics and historical performance analysis
  - Export statistics repository long-term
- ASH provides finer-grained session specific data
  - Catches variation in snapshot data
  - Stored in AWR repository
  - Used by ADDM

# Recommendations

- Most relevant data for analysis can be derived from the wait events

- *Always use EM and ADDM reports for performance health checks and analysis*

- ASH can be used for session-based analysis of variation

- Export AWR repository regularly to save all of the above

# For More Information

http://search.oracle.com

REAL APPLICATION CLUSTERS

or

otn.oracle.com/rac