

ORACLE

Exadata Database Machine : Maximum Availability Architecture (MAA)

Technical Presentation

Exadata and MAA Product Management

January 2025

Agenda

1

**Why focus on
Maximum
Availability?**

2

**What is
Maximum
Availability
Architecture?**

3

**Maximum
Availability
Architecture
features in
Exadata**

4

**Exadata
Lifecycle
Operations**

5

Summary

Why focus on Maximum Availability?



\$350K

—
Average Cost of downtime per hour

Source: Gartner, Data Center Knowledge, IT Process Institute, Forrester Research



\$10M

—
Average Cost of unplanned data center outage or disaster

Source: Gartner, Data Center Knowledge, IT Process Institute, Forrester Research

87 hours

—

Average Downtime per year

Source: Gartner, Data Center Knowledge, IT Process Institute, Forrester Research

91%

—

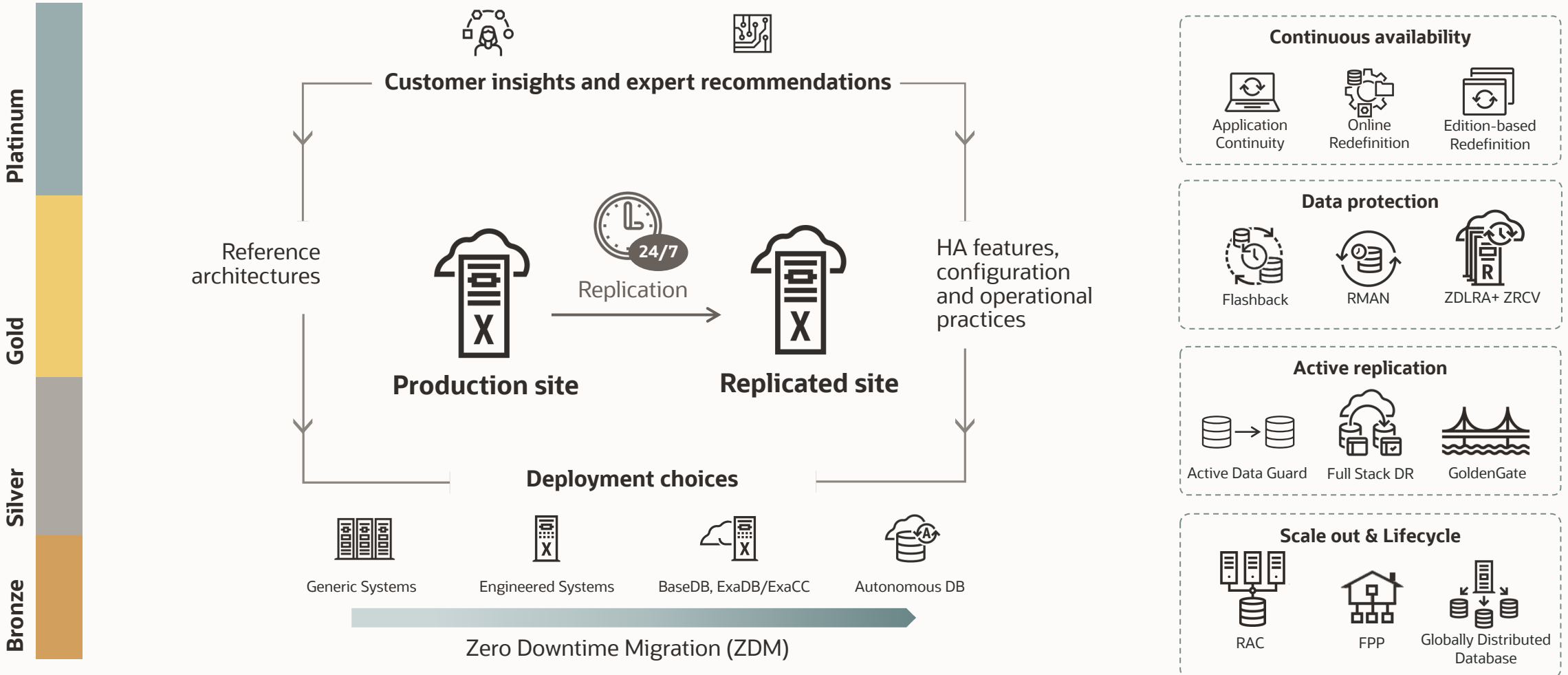
Percentage of companies that have experienced an unplanned data center outage in the last 24 months

Source: Gartner, Data Center Knowledge, IT Process Institute, Forrester Research

What is Maximum Availability Architecture?





Oracle Maximum Availability Architecture (MAA)

Standardized Reference Architectures for Never-Down Deployments



MAA Reference Architectures

Availability service levels

Bronze	Silver	Gold	Platinum
Dev, test, prod	Prod/departmental	Business critical	Mission critical
Single instance DB Restartable Backup/restore	Bronze + Database HA with RAC Application continuity	Silver + DB replication with Active Data Guard	Gold + GoldenGate Edition-Based Redefinition
			

All tiers possible with on-premises and cloud.



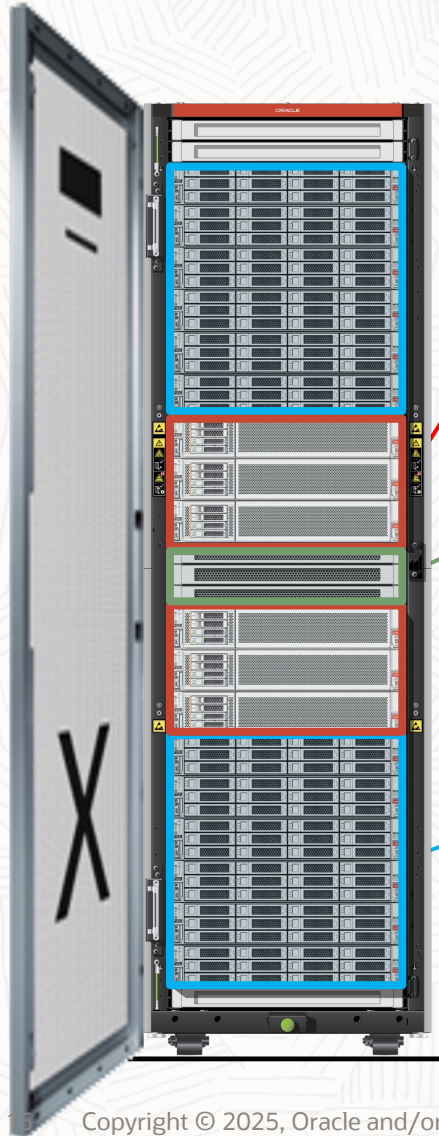
MAA Exadata Features

Hardware and Software Engineered together :

- ✓ Performance
- ✓ Manageability
- ✓ Availability



Oracle Exadata Database Machine : Built-in High Availability



● Redundant Database Servers

- Active-Active highly available clustered servers
- Hot-swappable power supplies, fans and flash cards
- Redundant power distribution units
- Integrated HA software/firmware stack

● Redundant Network

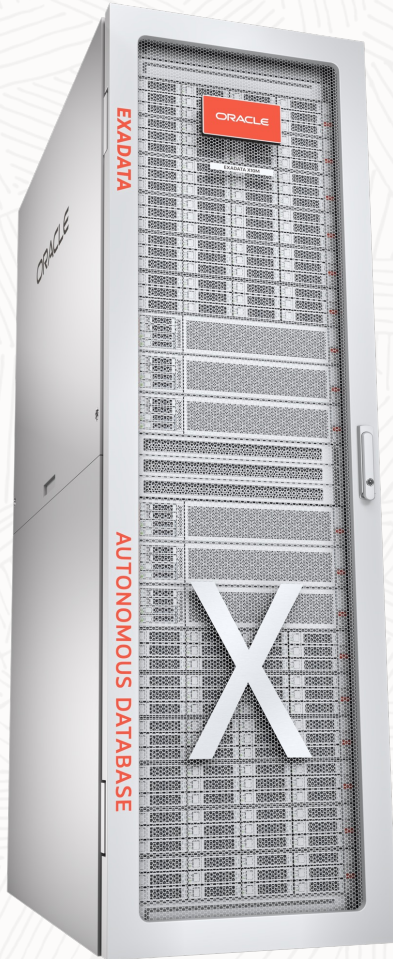
- Redundant 100Gb/s RoCE and switches
- Client access using HA bonded networks
- Integrated HA software/firmware stack

● Redundant Storage Grid

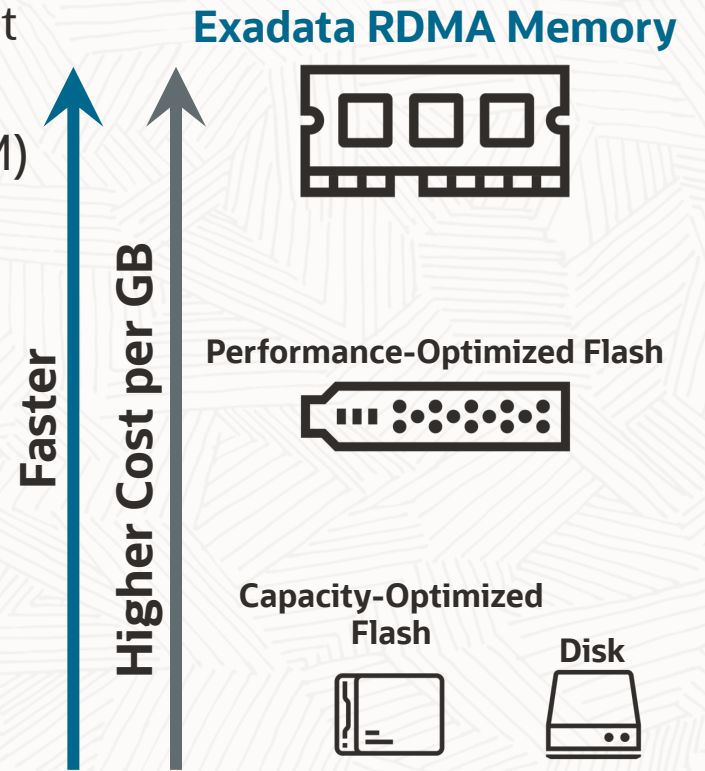
- Data mirrored across storage servers
- Hot-swappable power supplies, fans, M.2 drives and flash cards
- Redundant, non-blocking I/O paths
- Integrated HA software/firmware stack



Exadata : X11M



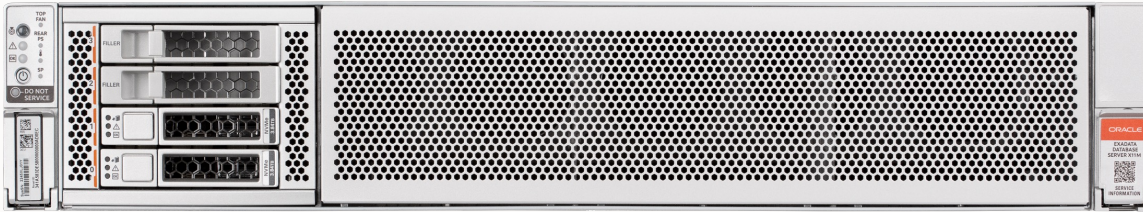
- 100 Gb active-active RDMA over Converged Ethernet (RoCE) private network
- 1.25TB low latency Exadata RDMA Memory (XRMEM) per Storage Server
- Data Acceleration reduces read latency to $<14\mu\text{s}$
- 3 storage tiers:
 - Exadata RDMA Memory
 - Performance-optimized Flash
 - Capacity-optimized Flash or Hard Disk
- Baremetal or KVM Based Virtualization



Exadata: The MAA Platform of Choice

Evolution: We Continue to Protect your Service Level from the Most Difficult HA Problems

X11M Database Server



Zero impact major Linux upgrades, e.g. OL8 in Exadata release 23.1

Zero impact security software upgrades including STIG compliance

MS (Management Server) alerting of key Database and Grid Infrastructure software incidents

14 microseconds to retrieve a database I/O from storage server XRMEM Cache

Human Error Prevention!

X11M Storage Server



Low I/O latency preservation during unplanned and planned outages

Tightly integrated hardware & software with auto repair of sick storage

Exadata X10M and X11M Extreme Flash storage server with both performance and capacity-optimized flash

MAA Best Practice Full Stack Compliance Checks with Exachk

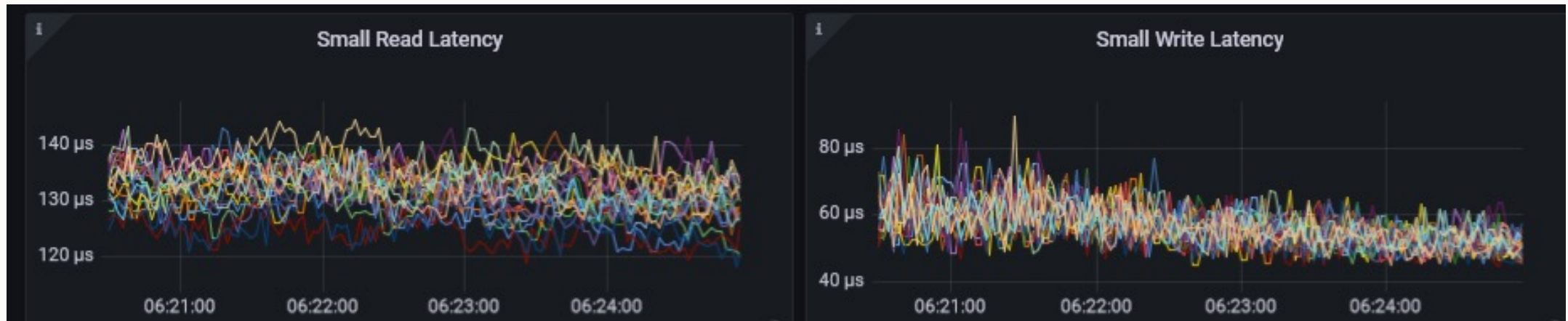


Exadata: The MAA Platform of Choice

Evolution: Metrics and More Made Easy

How do I really know what is going inside of Exadata?

- Performance data including Exadata metrics have been around since Exadata inception but they were sometimes difficult to consume and understand
- Enter *Real Time Insight* in Exadata release 22.1. Simply zoom into one of the dashboards to observe performance trends or shine a bright light on performance anomalies



Exadata : Built-in High Availability

- Automatic LED support for disk removal
- Redundancy Check during power down
- I/O error prevention with Exadata disk scrubbing / ASM corruption repair
- Failure Monitoring on database servers
- Reduced brownout for instance recovery
- I/O latency capping for reads and writes
- ILOM hang detection and repair
- Custom Diagnostic Package for Cell Alerts
- Updating database nodes with patchmgr
- Optimized and Faster Exadata Patching
- Blue OK-to-remove LED light notification
- Automated repair from controller cache failure
- Exadata HARD
- Auto online
- Priority rebalance support
- Cell-to-Cell Rebalance Preserves Flash Cache
- Exadata Elastic Configuration
- Drop hard disk for replacement
- Cell to Cell offload for Disk Repair
- Flash and Disk Life Cycle Management Alerts
- Redundancy protection on cell shutdown
- Fast network failure detection
- Efficient resilver rebalance after flash failure
- Fastest Redo Apply and Instance Recovery
- Exadata Smart Write Back
- I/O hang detection and repair
- I/O and Network Resource Management
- Health factor on predicatively failed disks
- Exachk full stack healthcheck with critical issues alerts
- Elimination of false positive drive failures
- EM failure reporting
- Active Active ROCE Network
- VLAN support and automation
- Drop BBU for Replacement
- Exadata Smart Flash Logging
- Cell-to-Cell Rebalance Data Accelerator Cache preservation
- Disk confinement
- Redundancy protection on cellsrv shutdown
- Corruption prevention with HARD support
- Smart Write Back Flash Cache persistence
- Auto disk management
- Cell I/O timeout threshold
- Automatic ASM mirror read on I/O error corruption
- Appliance mode support
- Cell Alert Summary



Lifecycle Management

Data Protection

Brownout

Quality Of Service and Performance



Lifecycle Management

Data Protection

Brownout

Quality Of Service and Performance



What are Data Corruptions

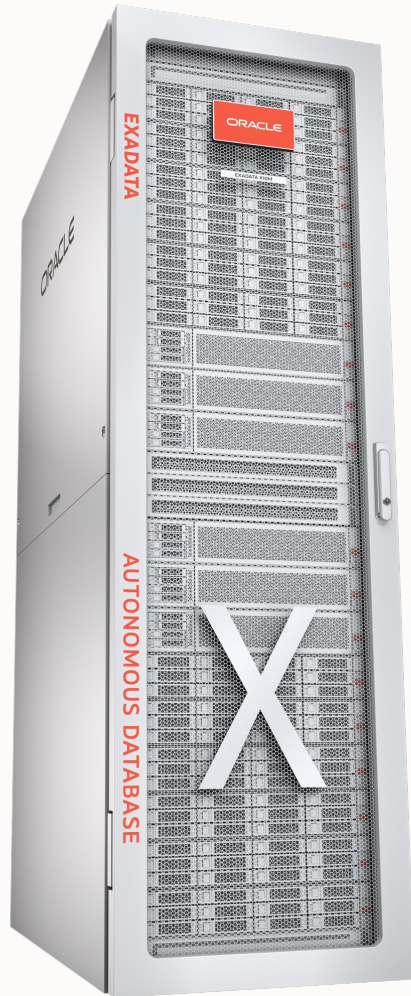
Data corruption refers to errors in **computer data** that occur during writing, reading, storage, transmission, or processing, which introduce unintended changes to the original data. Computer, transmission, and storage systems use a number of measures to provide end-to-end **data integrity**, or lack of errors.

- **Physical corruptions aka Media Corruption**
 - Checksum of the database block doesn't correspond with its contents
 - Header corrupt
 - Block Contains Zeros
 - ...
- **Logical corruptions**
 - Database block with correct checksum but logically inconsistent
 - Structure below the header is corrupt
 - Lost Write
 - Row locked by inexistent transaction
 - ...
- Often occurs silently



Exadata : Data Protection

Corruption Detection & Prevention



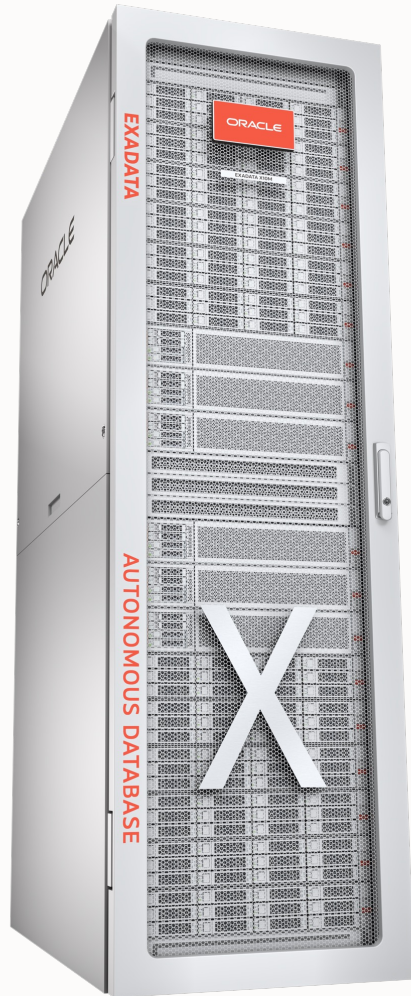
When a network packet in the I/O path between DB server and storage node is corrupted

- Storage cell prevents the write
- ASM retries by re-sending the packet
- ✓ Application never encounters corruptions



Exadata : Data Protection

Corruption Repair



If an application update in the database encounters corruption

- Database reads from the ASM mirror
- Repairs the corruption using the good copy
- ✓ This repair happens without impacting other database processes and application



Exadata : Data Protection

Storage Failures

On the storage cells, what happens if a drive is reported as but has not really failed?

- Automatic power cycle the drive / flash to avoid false positive drive failure



- When a storage failure occurs, redundancy is impacted
- Restoration of redundancy is prioritized in database-aware order to preserve data
- Order of Priority
 1. Control Files
 2. Online logs
 3. Archivelogs
 4. ASM SPfile
 5. Database SPfile
 6. TDE key store
 7. OCR
 8. Standby Redo Logs
 9. Wallet
 10. Datafiles

Exadata : Data Protection

Efficient Rebalance with Service Level Protection



ASM Power Limit: Intelligent and flexible rebalance power setting

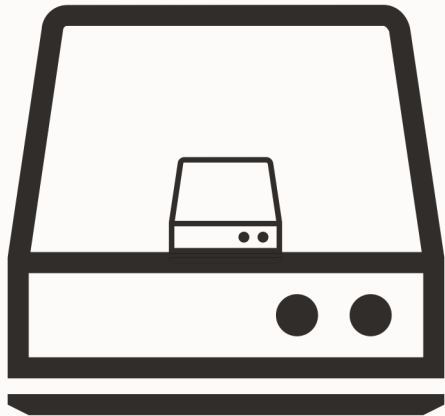
- Testing in MAA labs to find the best balance between redundancy restoration and service level protection
- MAA best practice `asm_power_limit`
 - Default 4 (total across clusters) at deployment time
 - Never set `asm_power_limit = 0`
 - Dynamically modifiable using – see table for recommended max
`alter diskgroup <diskgroup name> rebalance modify power <value>`

Recommended MAX <code>asm_power_limit</code>	
Oracle Database 23ai	Oracle Database 21c and earlier
96	64



Exadata : Data Protection

Exadata ASM configuration best practices

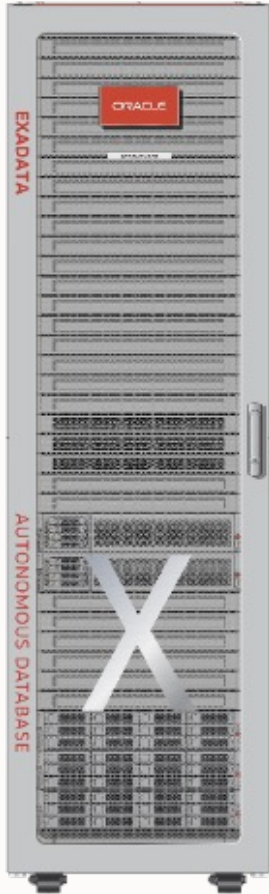


High Redundancy HIGHLY recommended

- Disks are constantly getting larger
 - During rolling software cell updating 2 copies remain
 - Double partner disk failures are rare but possible
- Particularly important for older systems with **aging** disks
- User data is stored in Primary Extent and 2 mirrored copies
- 5 failure groups required for Voting files and ASM metadata
 - OCR stored in primary and 2 mirrored copies

Exadata : Data Protection

Exadata ASM configuration best practices



High Redundancy requires at least 1 disk group with 5 failure groups

- An Eighth and Quarter Rack has 3 Storage Servers
- Only Storage Server 3 failure groups

Solution : **ASM quorum disk on Database Servers**

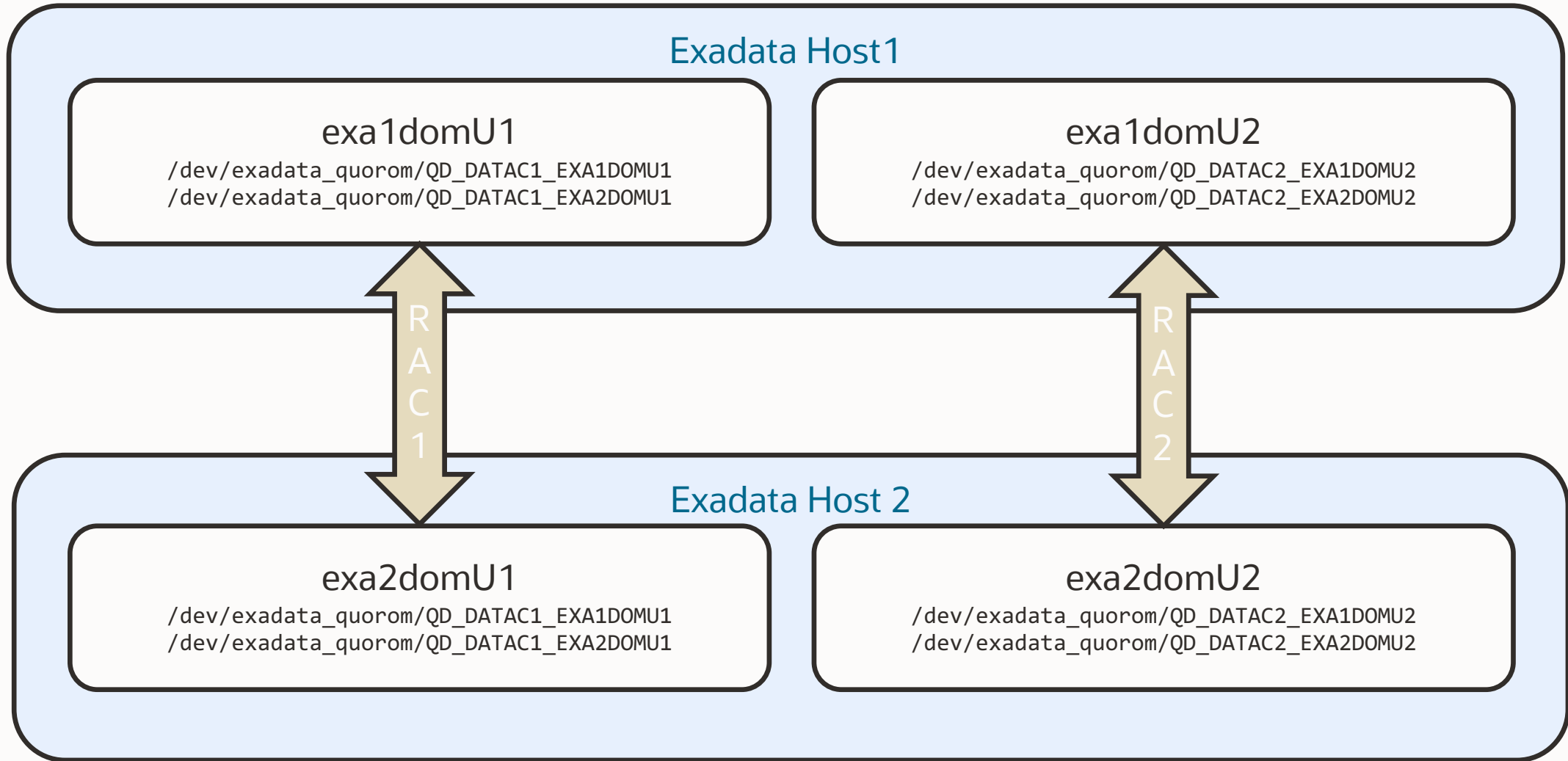
- Implemented automatically when deployed through OEDA
- Uses iSCSI based 'quorum failure groups'
- Managed with `quorumdiskmgr` (if needed)

High Redundancy HIGHLY recommended



Exadata : Data Protection

Exadata ASM configuration best practices Eighth & Quarter Rack High Redundancy



Exadata : Data Protection

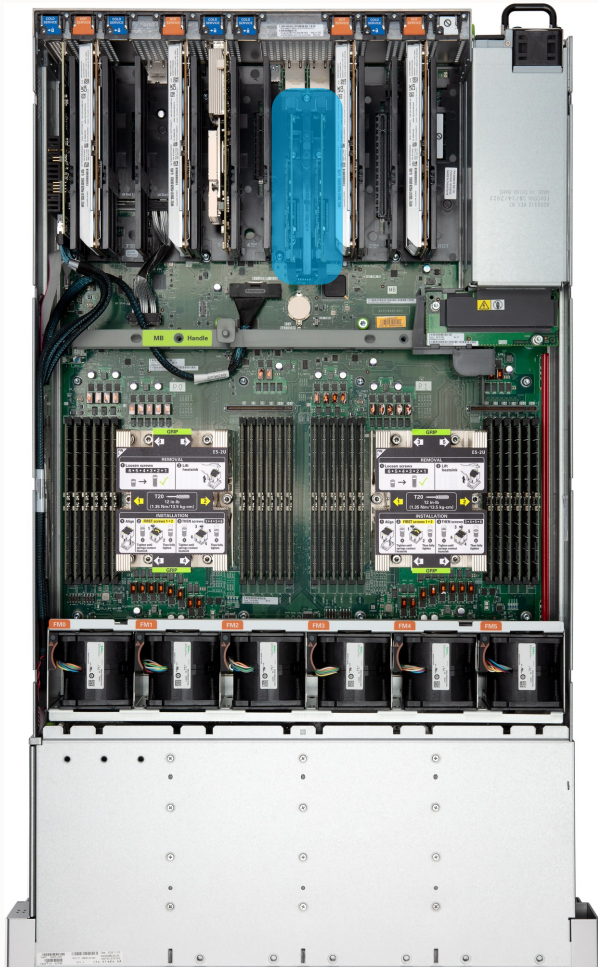
Storage Failures

- Exadata includes automated operations for disk maintenance when disks fail or have been proactively marked as problematic
- ASM automatically restores redundancy before balancing data on disk
 - Reduces window when some data may have reduced redundancy
- If a disk needs to be dropped manually, administrator can specify `MAINTAIN REDUNDANCY` to rebalance data before dropping the corresponding ASM disks
 - Preserves redundancy in addition to regular checks performed by `DROP FOR REPLACEMENT`



Exadata : Data Protection

M.2 Fast Failure Protection and Online Replacement (X7 and newer)



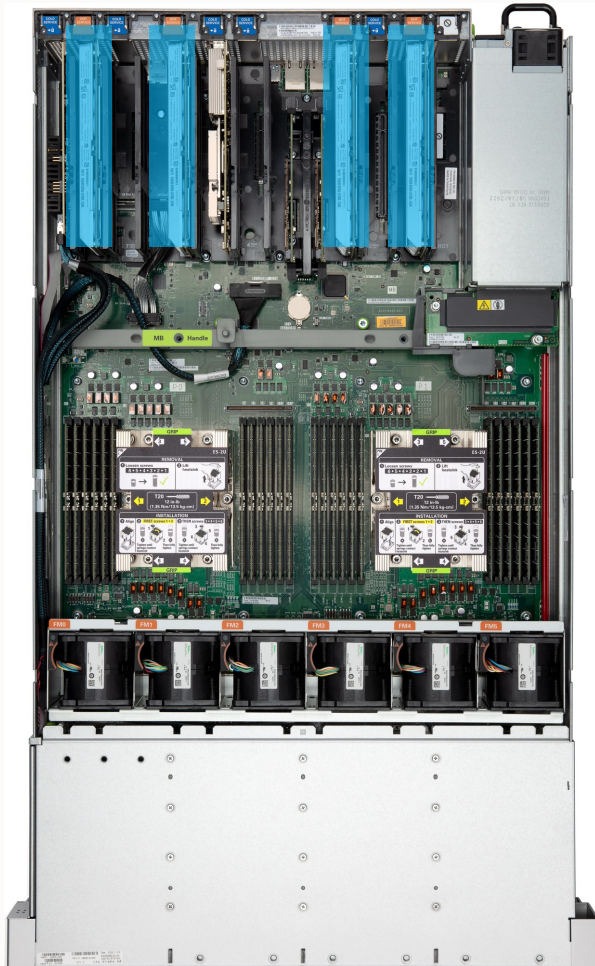
Two M.2 drives for OS and cell software

M.2 drives protected with Intel RSTe Raid

Can be replaced online so user data does not have to be taken offline

Exadata : Data Protection

Online Flash Replacement (X7 and newer)

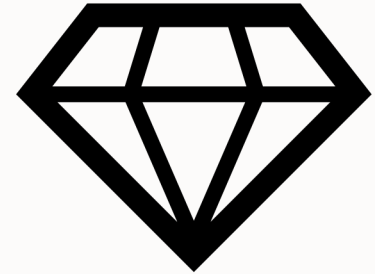


- Open chassis and replace online; no outage needed from storage server
- For failed drive replace when ready
- For online drive :
 - CellCLI> alter physicaldisk FLASH_2_2 drop for replacement;
- After replacement no customer interaction needed

Exadata : Data Protection

Hardware Assisted Resilient Data

- Exadata includes Hardware Assisted Resilient Data (HARD) checks to prevent corruption for specific file types:
 - Spfile
 - Controlfiles
 - Log files
 - Datafiles
 - Data Guard Broker Files
- When HARD check fails corrupted data not written
- Works transparently after enabling DB_BLOCK_CHECKSUM
 - Active during ASM Rebalance or ASM Resync



Exadata : Data Protection

Disk Scrubbing

- Inspects and repairs hard disks during idle time
 - Checks for bad sectors on the disks
 - Executed by Exadata Storage Software
 - If bad sectors are found storage requests mirror copy from ASM to perform repair
- Automatic and dynamic execution
 - Scheduled by default bi-weekly
 - When disks are idle (< 25% busy)
 - Automatically backs off when application needs I/O resources



Conclusion Data Protection

Rest assured Exadata has you covered :

- Corruption Detection, Prevention & Repair
- H.A.R.D.
- Scrubbing
- Online flash replacement
- M.2 Fast Failure Protection and Online Replacement
- High Redundancy
- Do not service LED
- Efficient Rebalance
- Automatic Power cycle of (potentially sick) flash/ drives



Credit : David Clode
https://unsplash.com/photos/Yg_sNKOixvY



Lifecycle Management

Data Protection



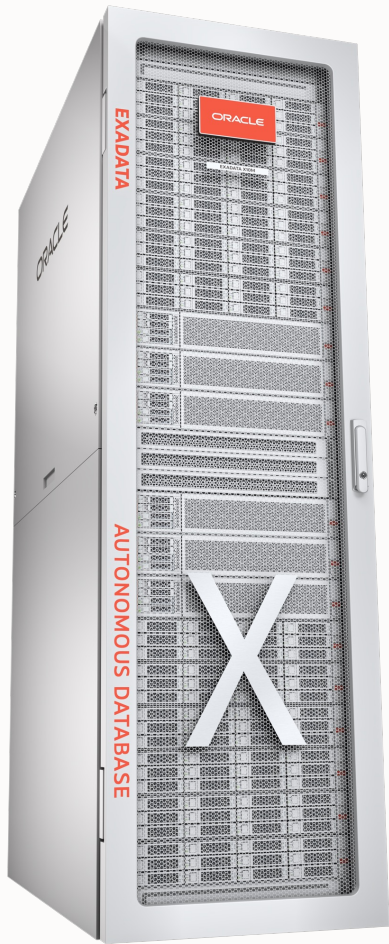
Brown out

Quality Of Service and Performance



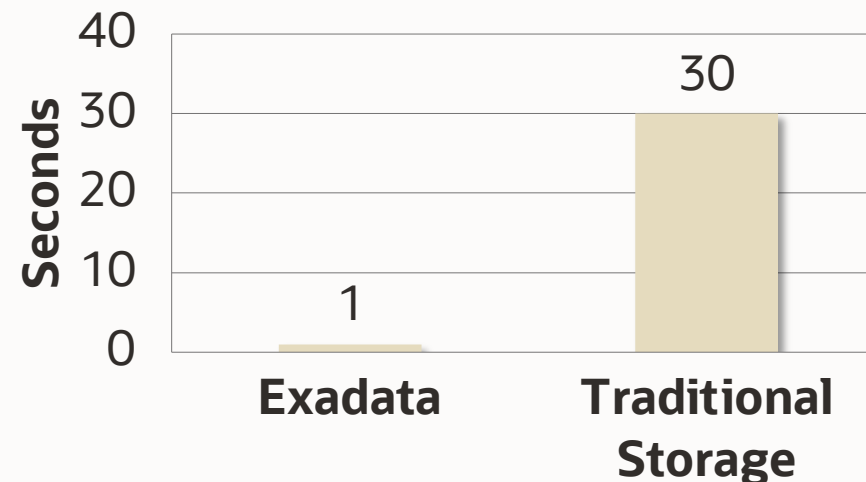
Exadata : Quality of Service & Performance

I/O Latency Capping



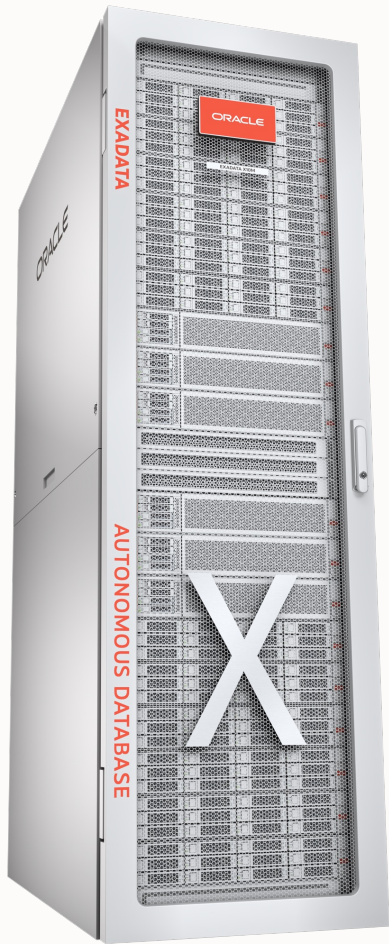
- High I/O latency can have detrimental performance impact
- Exadata detects high I/O latency and redirects reads and writes to other devices
 - High latency read I/O redirected to partner cell
 - High latency write I/O cancelled and temporarily written to flash on same cell

LGWR Delay after Hung IO



Exadata : Quality of Service & Performance

Storage Server Disk Confinement



- Exadata constantly monitors disk performance and health
- Poor performance is often a precursor to disk failure
- Disks identified with poor performance are confined and I/O directed to alternative mirror
- Storage Server automatically runs disk health check
- If the disk is deemed healthy
 - Disk is returned to service and RESYNCRonized
- If the disk is deemed unhealthy
 - Disk is dropped, data rebalanced to maintain redundancy and blue service LED is lit
 - Disk can then be replaced



Exadata : Quality of Service & Performance

Smart Storage with I/O Resource Manager (IORM)

- IORM configures and manages Storage Server I/O related resources when contention occurs
- I/O tagged and prioritized based on IORM Plan for Database (CDB/PDB/Non-CDB) or Cluster
- Tag includes
 - Database/PDB/Cluster name
 - Purpose
 - Priority
- Useful in mixed and consolidated workload environments
- Can be combined with Database Resource Manager



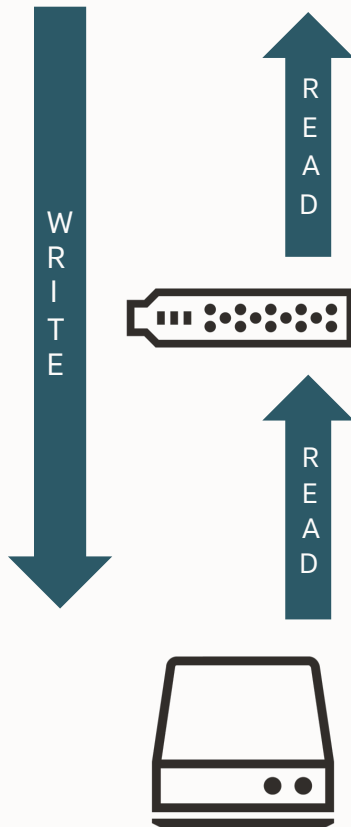
```
CellCLI> ALTER IORMPLAN - dbplan=((name=prod, share=16), -  
(name=dw, share=4), - (name=prod_test, share=2), -  
(name=DEFAULT, share=1))
```



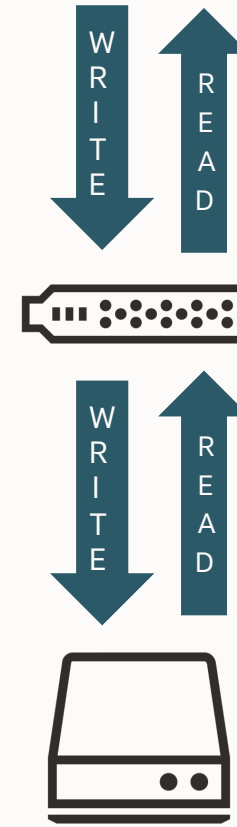
Exadata : Quality of Service & Performance

Flash Cache : Write Back or Write Through

Write-Through



Write-Back



Application I/O rarely hits the hard disks or Capacity-Optimized Flash

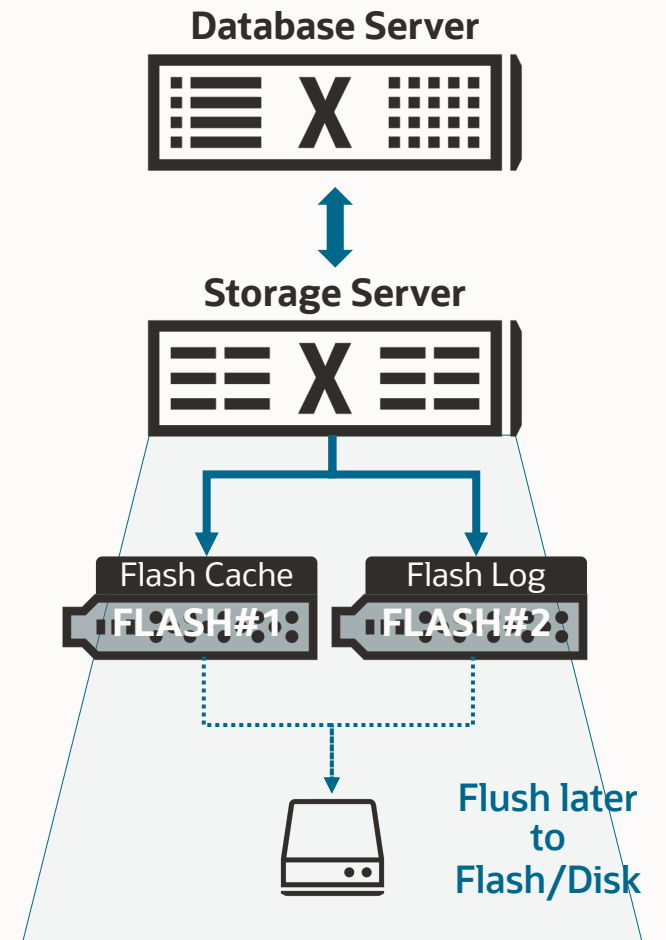
MAA Best Practice



Exadata : Quality of Service & Performance

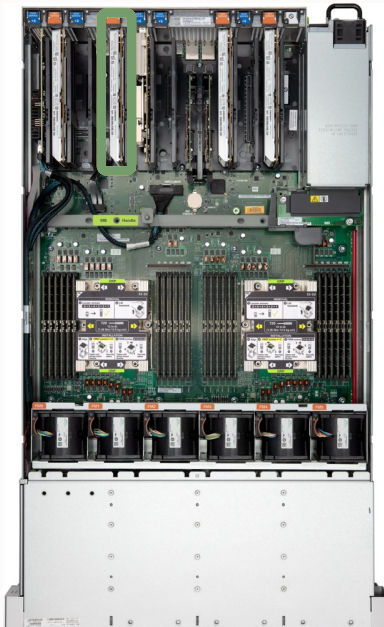
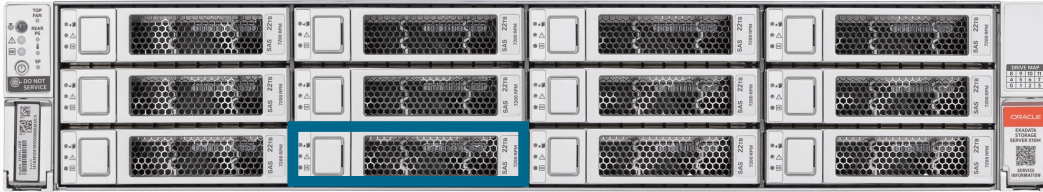
Smart Flash Log

- Eliminates write latency outliers
 - Redo log write latency is critical for OLTP database performance
 - Writes to Flash Cache and Flash Log on different flash devices simultaneously
 - Fastest device acknowledges write
- Eliminates storage as log write bottleneck
 - Online and standby redo logs automatically and transparently cached in write-back Smart Flash Cache
 - Increases log write throughput by writing to flash instead of disk
 - Benefits workloads that read the online redo logs such as GoldenGate
 - Beneficial when multiple concurrent workloads require hard disk I/O bandwidth (eg backups)
- Asynchronous flush to capacity-optimized flash or HDD



Exadata : Quality of Service & Performance

Smart Flash and Hard Disk Replacement



- After flash or hard disk replacement, a “health factor” is set on the affected hard disks.
- While the health factor is on, reads are satisfied from a **healthy partner cell** and Exadata software continues **warming up the flash cache** on the cell that had its storage replaced.
- When the flash cache is sufficiently warmed up, the health factor status is removed.
- This feature enables **consistent, low I/O latency** after storage replacement that in turn **maintains application service levels**.

Exadata : Quality of Service & Performance

SLA's maintained during planned maintenance or unplanned maintenance

- Exadata flash cache state preserved during ASM rebalance operations. One practical example is the resync that occurs during cell software rolling updates.
- Intelligent routing of I/O requests to cell providing the best service after flash and disk failure and repair
- Applicable to both unplanned outages and planned maintenance

Performance is Time
Time is Money



Exadata : Quality of Service & Performance

Database Tier I/O Cancel


Database Tier



Database Tier I/O Latency Capping ✓

I/Os are Pumping

Storage Tier



- Slow I/O ? Cell I/O Latency Capping ✓
- Hung I/O ? I/O Hang detection / repair ✓
- Sick disk ? Disk confinement ✓
- Undiscovered hardware / software issue?



Exadata ASM Reserved Space for Rebalance

- ASM requires space to allow for rebalancing of data in the event of a failure
 - Ensures rebalance is successful
 - Restores redundancy
 - Space to ensure rebalance is successful is not reserved
 - Reports ORA-15041 if there is not enough space to complete rebalance

REQUIRED_MIRROR_FREE_MB

- Depends on the number of failure groups and ASM version
Applies to any disk group and any redundancy (HIGH or NORMAL)
- Same for all media types and hardware generations*

Grid Infrastructure Version	Number of Failure Groups	Required % Free of Disk Group Capacity
12.1.0	Any	15
12.2, 18.1+	less than 5	15
12.2, 18.1+	5 or more	9

* Exadata X10M and newer Extreme Flash has hardware-specific requirements

Number of Failure Groups (8 ASM disks / FG)	Redundancy	Required % Free of Disk Group Capacity to Successfully Rebalance after a single physical disk failure
less than 5	NORMAL	15%
less than 5	HIGH	29%
5 or more	NORMAL	9%
5 or more	HIGH	11%

- X10M and newer EF cells have four physical flash disks with two ASM disks per physical flash disk. Therefore, a flash card failure will result in two ASM disks being dropped.
- GI/ASM 19c and newer with patch 34281503



Exadata Smart Rebalance

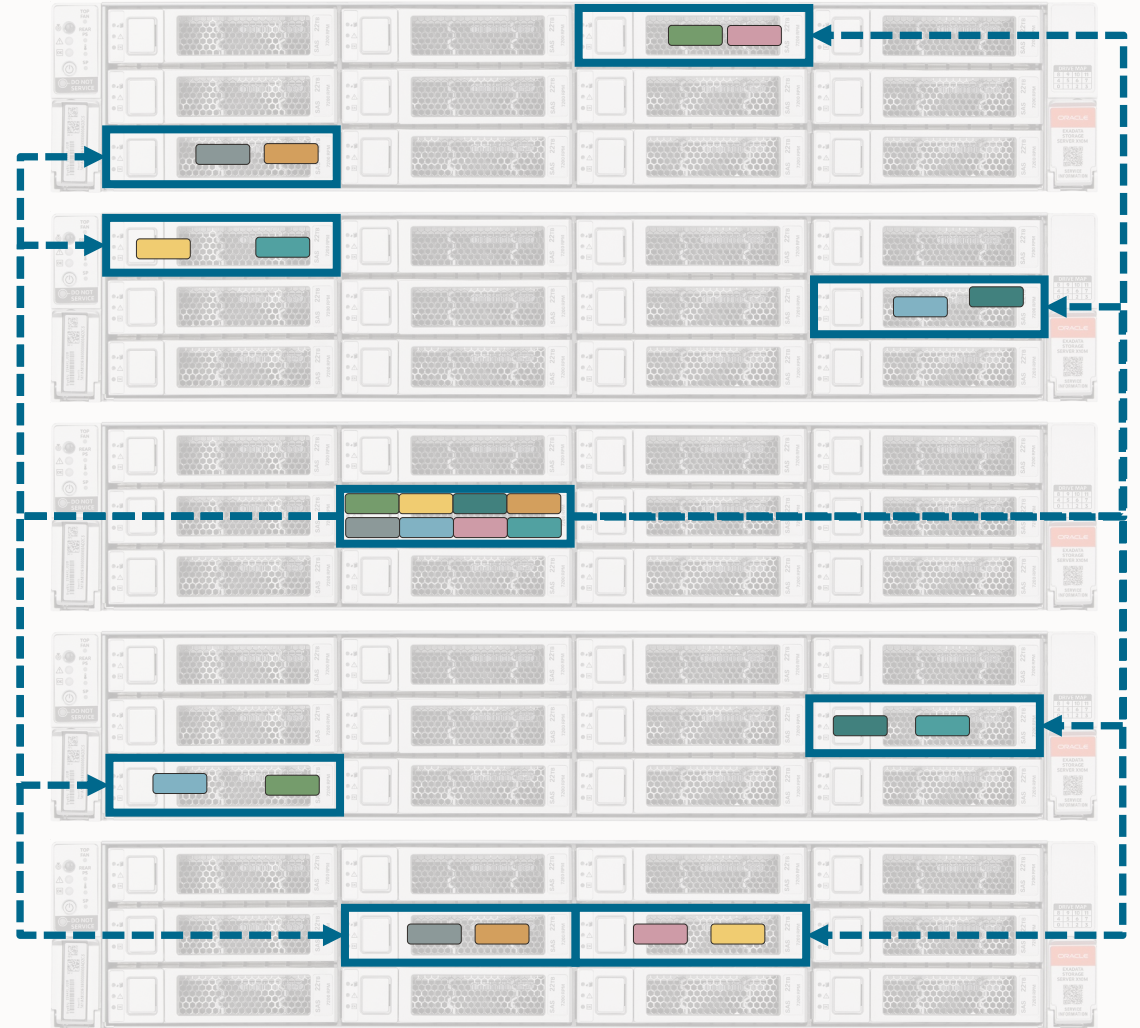
- Smart Rebalance affects High Redundancy disk groups when a failure occurs
 - If disk group has required free space
 - Data is rebalanced and redundancy restored
 - If disk group DOES NOT have required free space
 - Disk is offlined and rebalance deferred
 - Disk is re-mirrored efficiently from partner disks once replaced
- Reduces data movement and extra I/O at failure time if more capacity is required for database storage

Smart Rebalance is a safety-net. MAA strongly recommends maintaining sufficient free space



ASM Disk Partnering Concept

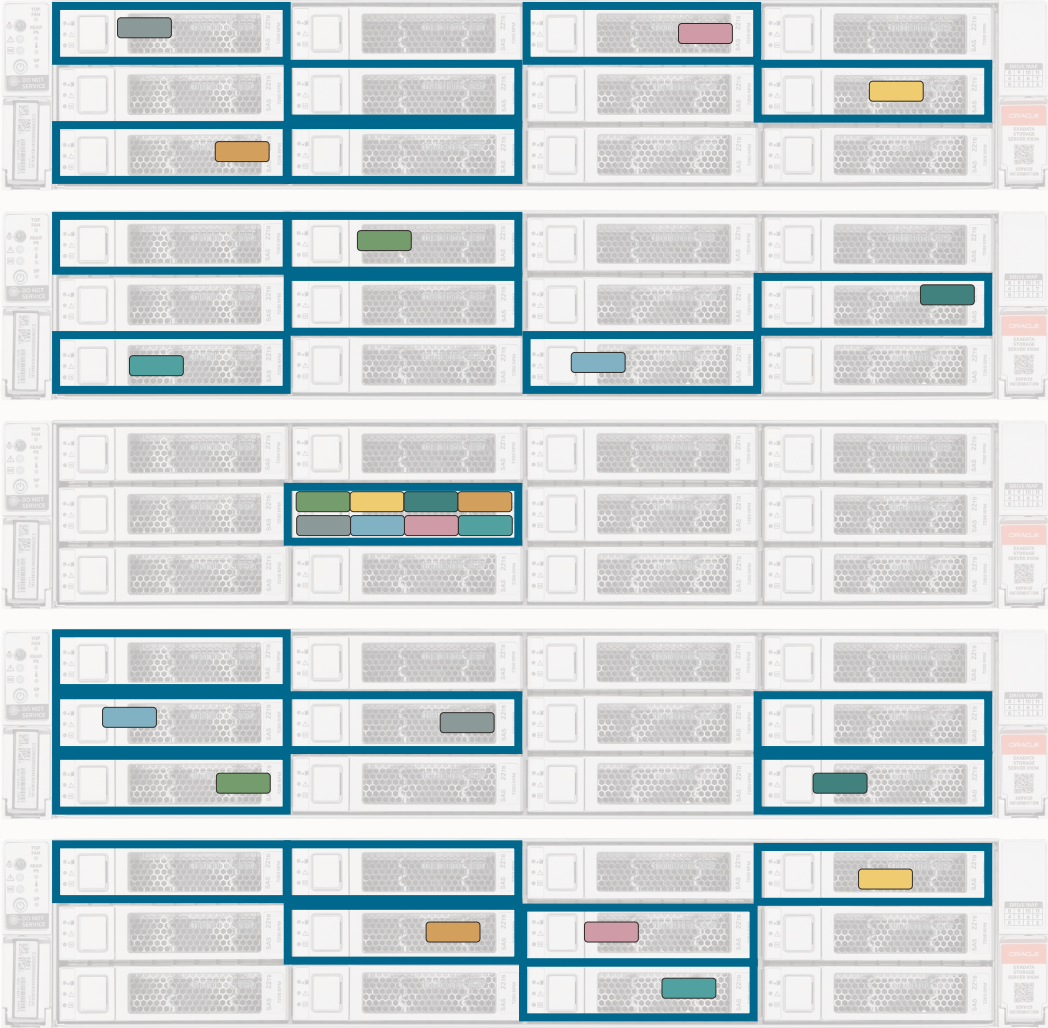
- ASM utilizes disk partnerships to choose disks for placing extents and their mirror copies
- Each disk partners with 8 other disks
 - For less than 5 cells, all partners are from 2 cells
 - For 5 or more cells, all partners are from 4 cells
- The Primary Extent is then mirrored
 - to two of these partners for High Redundancy
 - to one of these partners for Normal Redundancy
- Read IO provided by 8 disks
 - Used by rebalance, rebuild, resync, resilver, disk/flash warmup operations



ASM 23ai Increases Number of Disk Partners

- Each disk partners with 24 other disks on four cells
- Read IO provided by 24 disks
 - Benefits rebalance, rebuild, resync, resilver, disk/flash warmup operations
- Automatically managed by ASM
 - New partnering scheme not applied during upgrade
 - Partners updated by following operations
 - ADD DISK
 - ADD FAILGROUP (add cell)
 - REBALANCE

Results in up to 3x
faster redundancy restoration



What about other storage configurations?

- Different Storage Server configurations utilize different partnership values

Number of cells	Storage Server Type	Number of disks per cell	Number of partners (pre-23ai)	Number of partners (23ai)
3	1/8 th Rack High Capacity	6	8	12
3 or 4	High Capacity	12	8	12
5 or more	High Capacity	12	8	24
3 or 4	Extreme Flash	8	8	8
5 or more	Extreme Flash	8	8	16

- Benefits following operations
 - REBUILD – disk failure
 - RESYNC – cell patching
 - RESILVER – flash card failure
 - Disk/Flash WARMUP

Note – when adding a 5th cell to a configuration, rebalance will run longer as the increased number of partners is applied



Exadata : Quality of Service & Performance

Capacity Planning : Memory Configuration

Memory swapping can cause performance and stability issues

Correct memory configuration avoids :

- Swapping
- Instability



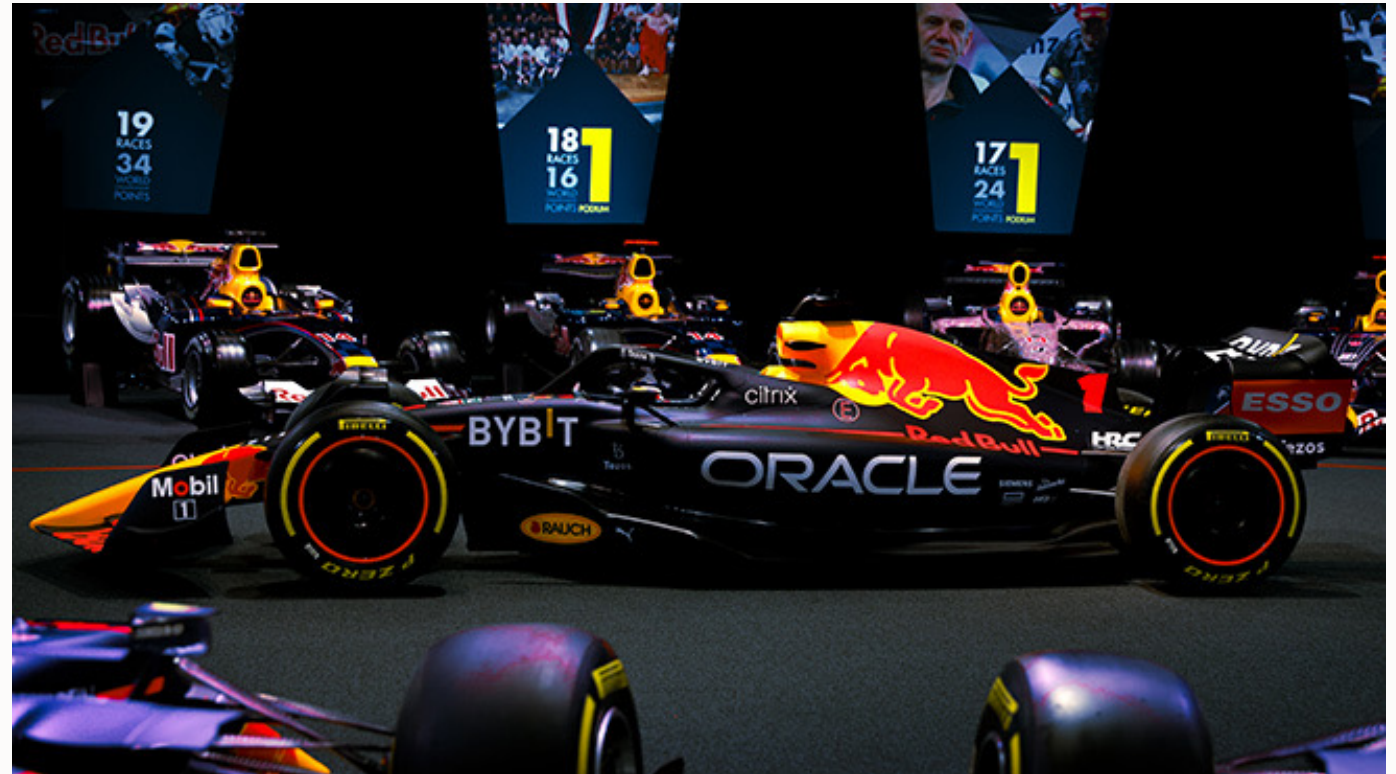
Credit : Kathy
<https://unsplash.com/photos/R7nSPG8edVI>



Exadata : Quality of Service & Performance

Exadata built for speed

- Smart Scan
- Smart Flash Cache
- Storage Index
 - “The fastest I/O operation is the one that you don’t need to do”
- Hybrid Columnar Compression
- In-Memory Columnar Format
- RDMA
- Real-Time Insight



Exadata : Quality of Service & Performance

RDMA Network Fabric

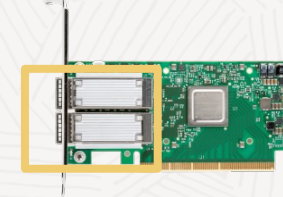
2 Active – Active ports in every RDMA Network Fabric Adapter

2 RDMA Network Fabric Switches in every Exadata single rack

22 Ports per switch used for internal cluster network, cabled ensuring no single point of failure exists

- Only to be used for Exadata purposes
- Settings on switch level not to be changed
- ZFS systems recommended to be connected through Top Of Rack (ToR) switches, for scalability and flexibility reasons

RDMA Network Fabric Adapter



RDMA Network Fabric Switch



Exadata : Quality of Service & Performance

Automatic Workload Prioritization

RDMA Fabric implements automatic Quality of Service (QoS)

Separate QoS lanes for specific traffic

- Critical I/O – LGWR
- Disk reads
- Disk writes



RoCE Network Resilience

Exadata RoCE IPs need to be highly available

- Each server has a dual-port RoCE NIC with each port connected to a different Leaf Switch
- Automatically failed over if a switch port is “down”

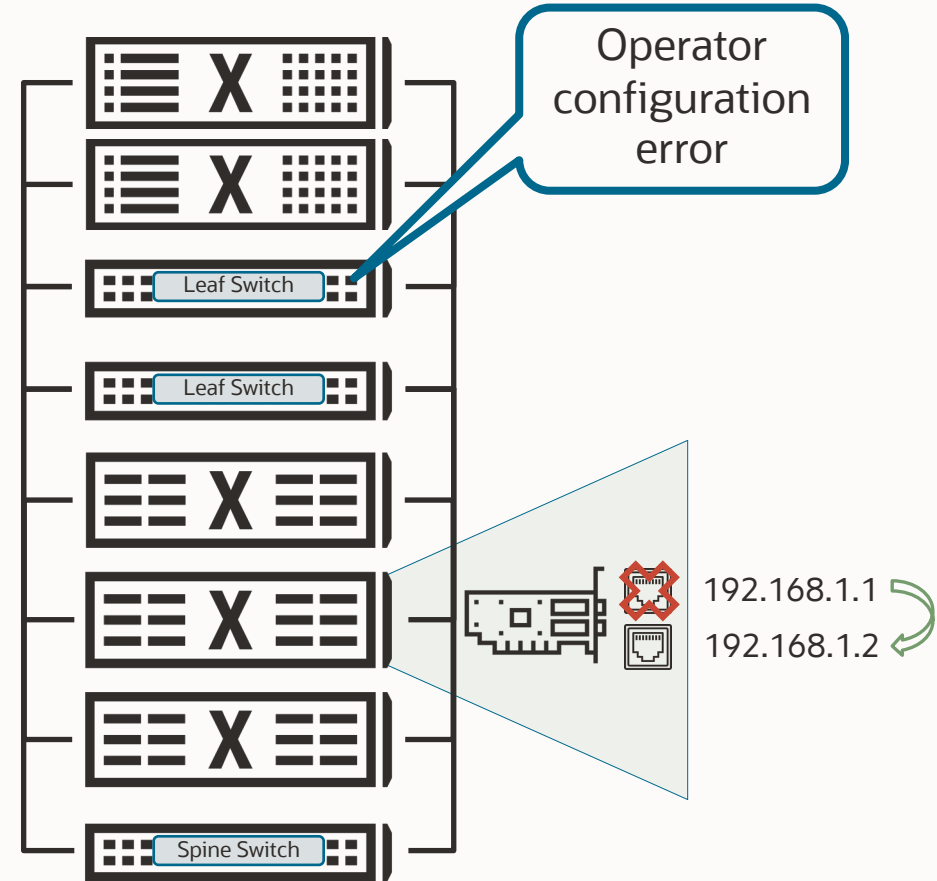
Unhealthy switches or network may leave ports “up” but network traffic stalled and unable to flow

- Switch misconfiguration
- Excessive pause frames

Network traffic stalls may result in database instability or outages

The **ExaPortMon** process runs on the host and monitors the live traffic of both RoCE ports

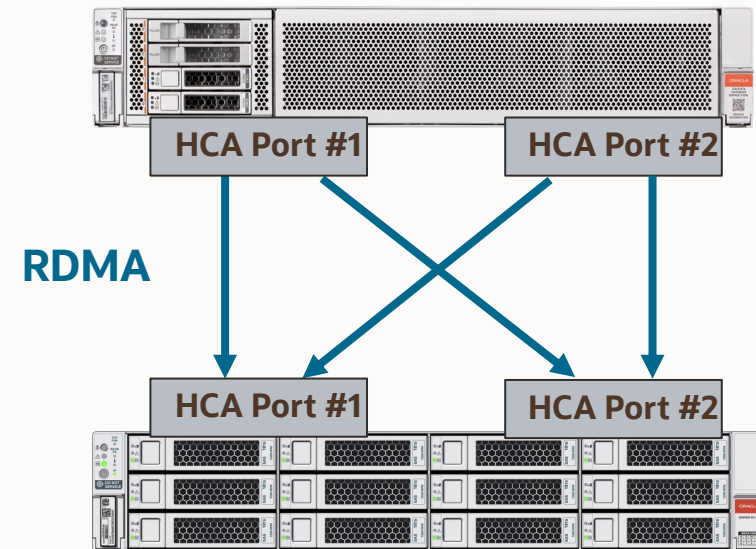
- Migrates IP to operational port if stall detected
- Returns IP to original port when upstream issue is resolved



Exadata : Quality of Service & Performance

Instant Failure Detection (IFD)

- Traditional systems use software to check availability
 - May cause performance issues under high load
 - Rely on TCP timeouts
- Exadata uses RDMA to check server availability
 - Instant Failure Detection
 - Utilizes 4 RDMA paths between for redundancy
 - Database ↔ Storage Servers
 - Database ↔ Database Servers
- If all four paths are unavailable after a short period the server is evicted



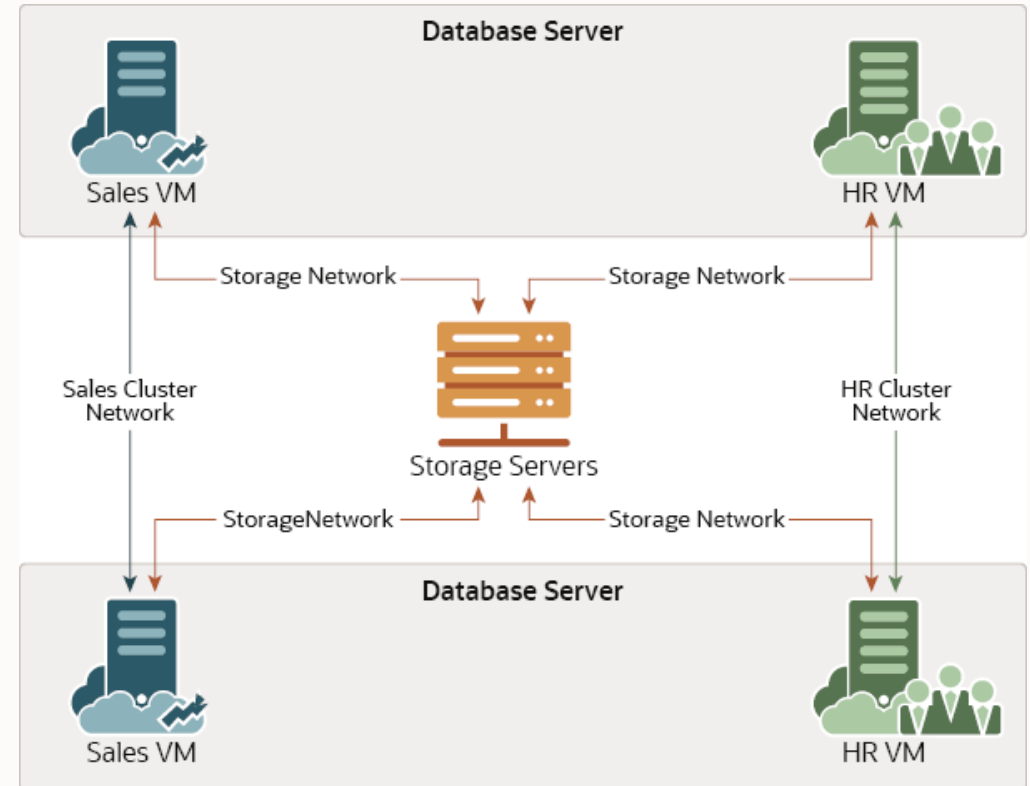
Sub-second notification vs up to 1-minute timeout on non-Exadata platforms



Exadata Secure RDMA Fabric Isolation for RoCE

Exadata **Secure Fabric** for RoCE systems implements network isolation for Virtual Machines while allowing access to common Exadata Storage Servers

- Each VM cluster is assigned a private network
- VM clusters cannot communicate with each other
- All VMs can communicate to the shared storage infrastructure
- Security cannot be bypassed
 - Enforcement done by the network card on every packet
 - Rules programmed by hypervisor automatically



Conclusion Quality of Service



Credit : Towfiq barbhuiya
<https://unsplash.com/photos/OZUoBTLw3y4>

- Cell side I/O Latency Capping
- Cell disk confinement
- Smart Storage with I/O Resource Manager IORM
- Smart Flash Logging
- Smart Flash Log Write-Back
- Smart Flash replacement
- Exadata RDMA Memory Data Accelerator
- RDMA Network Fabric
- RDMA QoS
- Instant Failure Detection
- Exadata Secure RDMA Fabric Isolation

Lifecycle Management

Data Protection

Brownout

Quality Of Service and Performance



Exadata : Brownout

Blackout vs Brownout

Blackout

- **Complete** service level interruption

Brownout

- **Significant** service level degradation

Lost productivity & Lost revenue

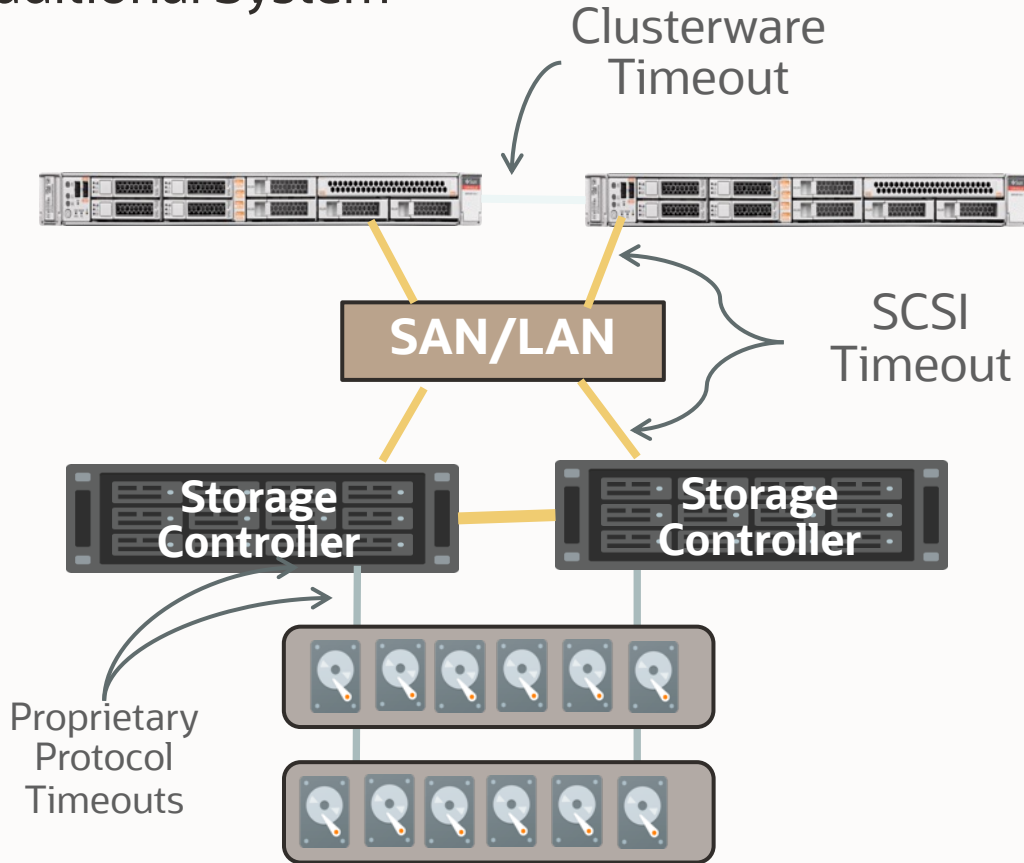
Systems are complex and an issue one layer can cascade to other layers

Our Engineered Systems and MAA best practices are designed and tuned to tackle this

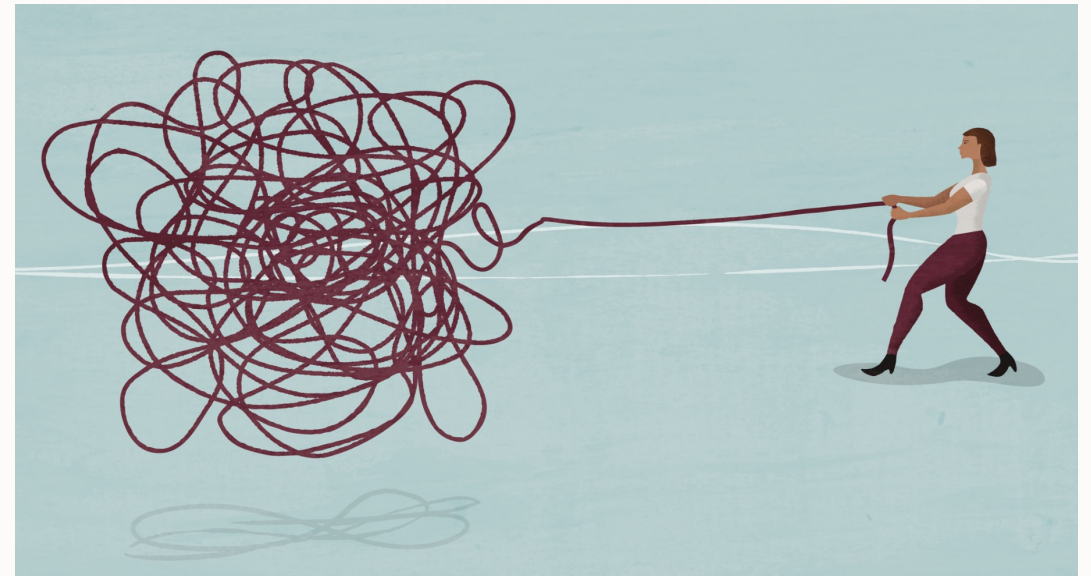


Exadata : Brownout

Traditional System



- Each layer has its own failure detection and timeouts
- Usually fault detection times additive
eg upon storage controller crash it takes 2 SCSI time-outs for db server to detect this failure



Cell Controller Cache Failure Handling

Automated Data Loss Prevention

Failed cache controllers can be complicated on custom built systems and earlier Exadata systems

Before Exadata 21.2, a user had to recover from a failed controller cache with manual steps:

- Answering cryptic question on the console about how to proceed
- Ensuring grid disks were force dropped before the controller was replaced



Cell Controller Cache Failure Handling

Automated Data Loss Prevention

Using Exadata 21.2 and higher, repair from controller cache failure is handled automatically by doing the following

- Detecting the problem before cell services start post crash
- Disable access to the grid disks
- Recover the failed disks



Credit :Connor McSheffrey
<https://unsplash.com/photos/MlspM6Hlit8>

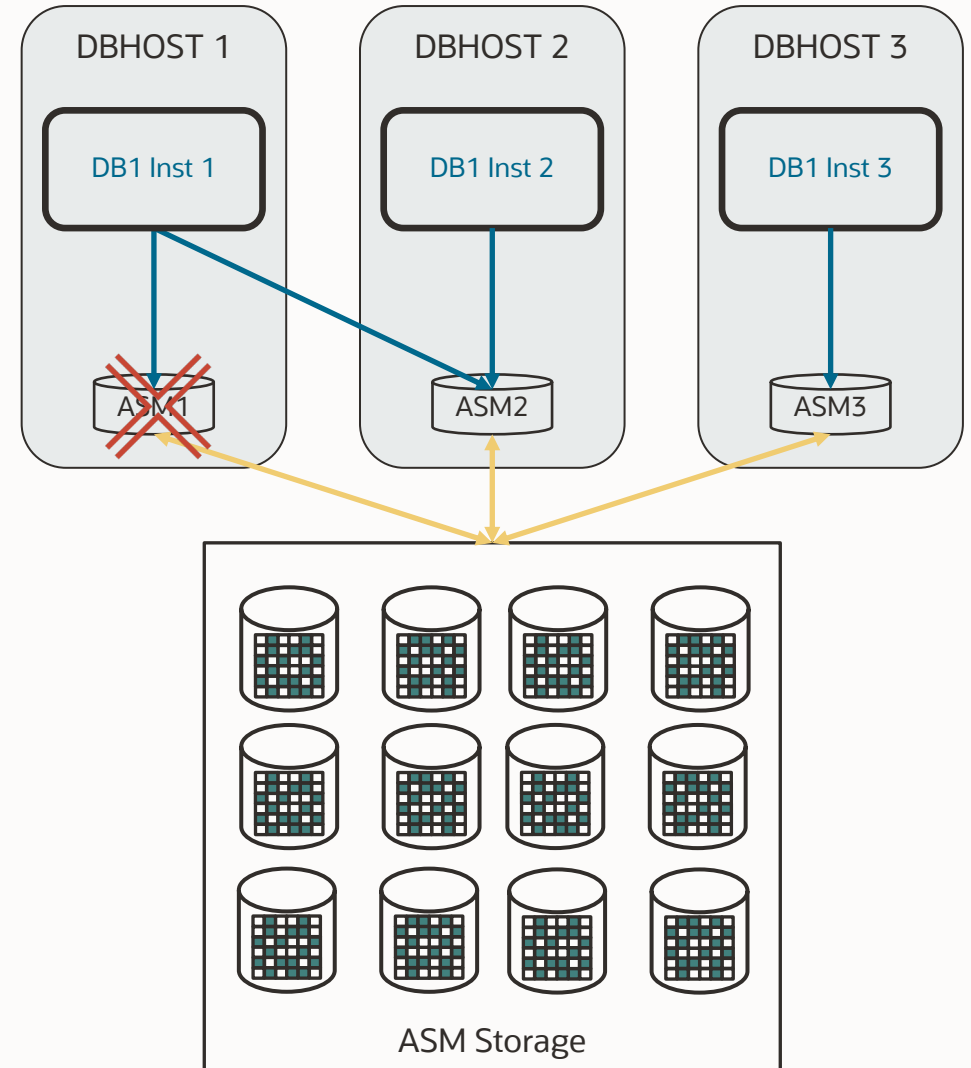


Exadata : Brownout

Brownouts & Blackouts : Flex ASM

Oracle Flex ASM enables Oracle ASM instances to run on a separate physical server from the database servers.

- Enables continuous RDBMS ↔ ASM communication
- After ASM instance crash no need for a service failover
- Completely transparent to the application with no service level impact
- On Exadata Cardinality is set to ALL

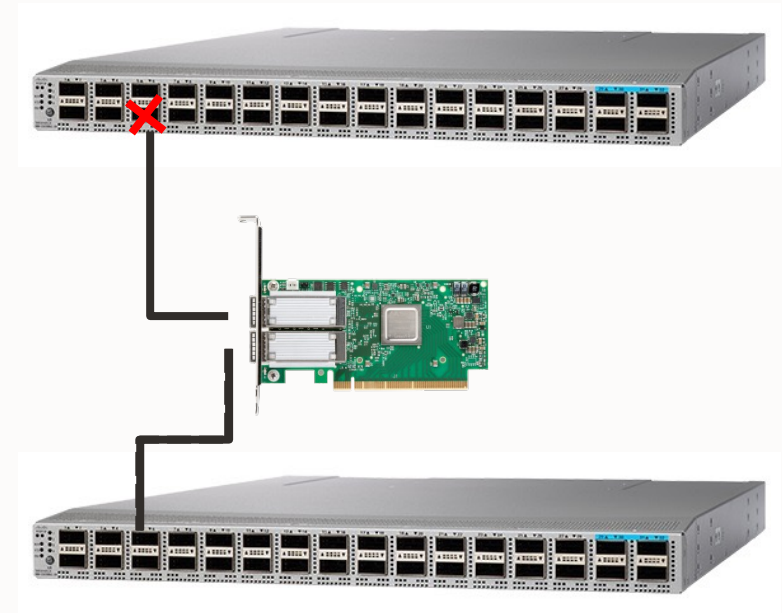


Exadata : Brownout

Brownouts Reduction for Client Network Port Failure

Brownout associated with active/passive client access network port failure is extremely low.

LACP “Active / Active” can also be configured but needs changes on network infrastructure



Exadata : Brownout

Smart Handshake for Storage Server Shutdown

- When storage server is shutdown the diskmon process in the Grid Infrastructure on the database server is notified
- No blackout when storage tier is shutdown for maintenance

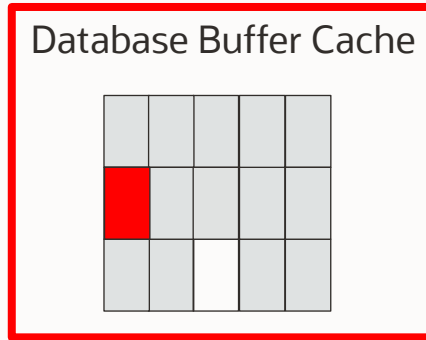


Exadata : Brownout

Smart OLTP Caching

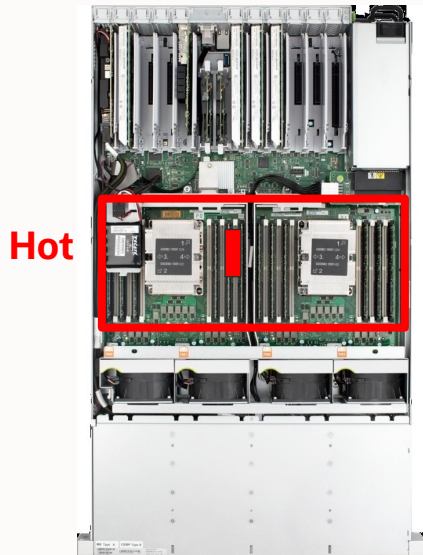
Quick review of Exadata data access tiers first...

- 1. Data read into buffer cache
- 2. DBWR evicts a buffer to free up space in buffer cache



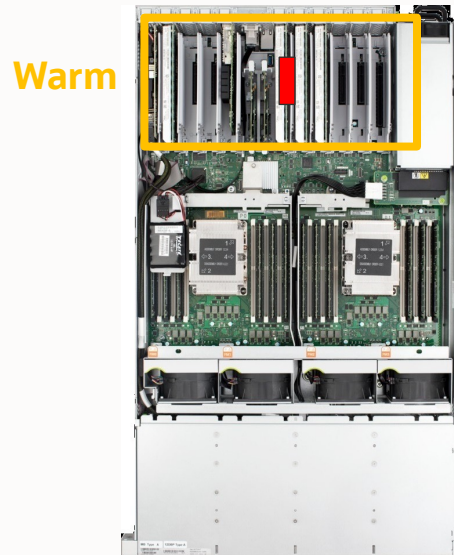
Sizzling

- 2. Cell with primary mirror populated in super low latency Data Accelerator



Hot

- 2. Cell with secondary mirror populated in low latency flash cache



Warm

- 2. Cell with tertiary mirror located on high latency hard disk throughout



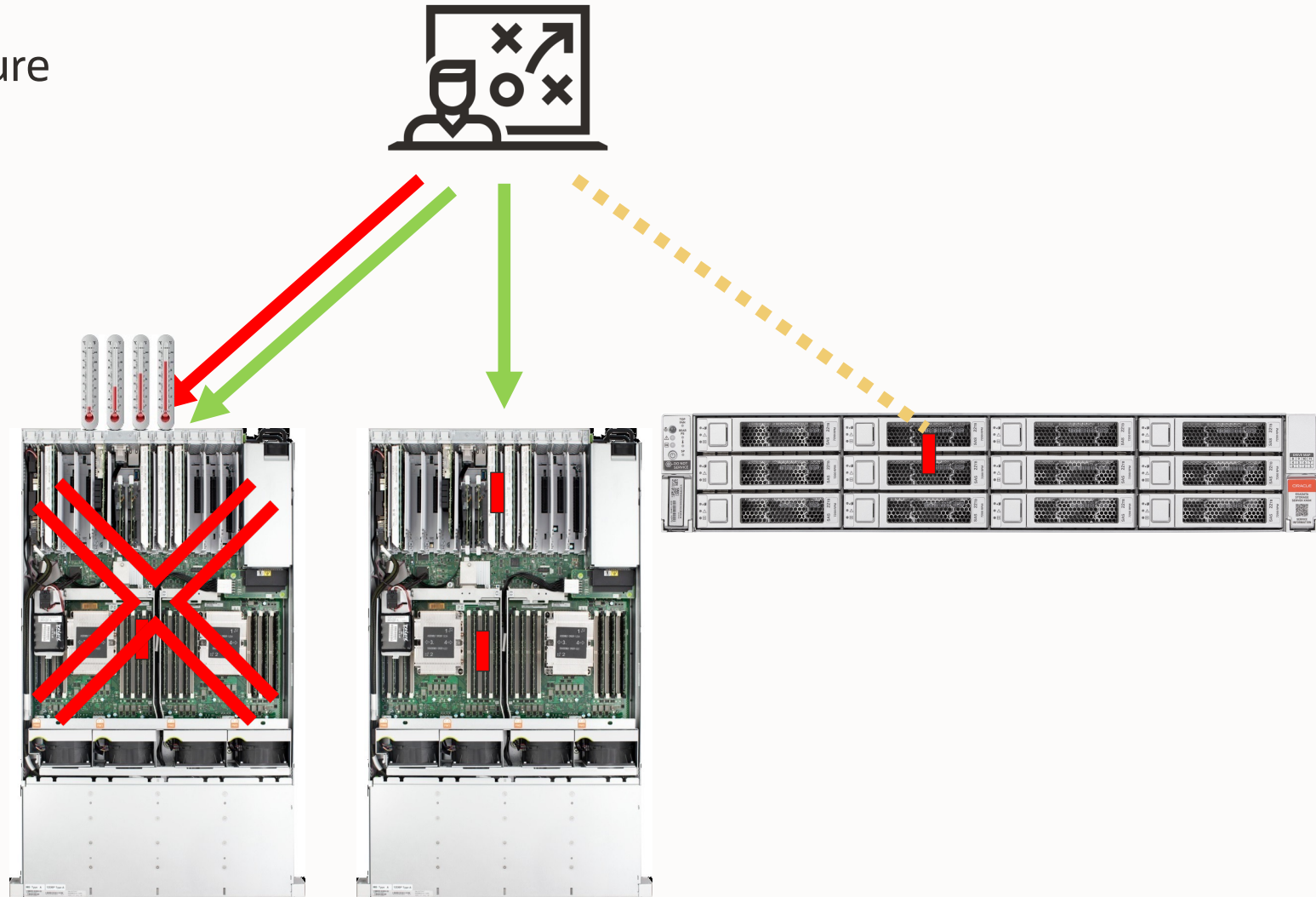
Cold



Exadata : Brownout

Smart OLTP Caching – Storage Failure

- Application reading data from primary mirror
- Storage failure on cell containing primary mirror
- Retrieve data from secondary mirror on flash with low latency and populate super low latency Data Accelerator
- Tertiary mirror continues to provide protection when Murphy strikes
- After repair of storage failure and flash cache warm up, return to primary copy



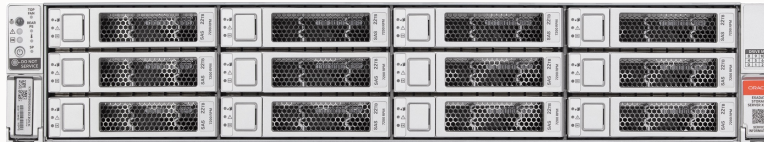
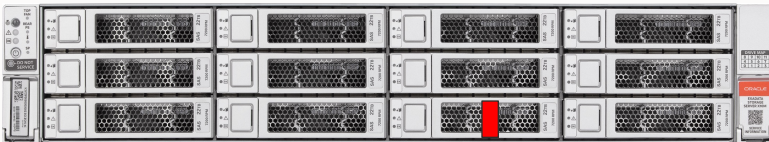
ASM rebalance, resync, and resilver always preserve flash cache state when moving extents



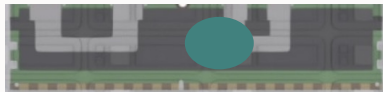
Exadata : Brownout

Cell-to-Cell Rebalance Preserves Data Accelerator Population

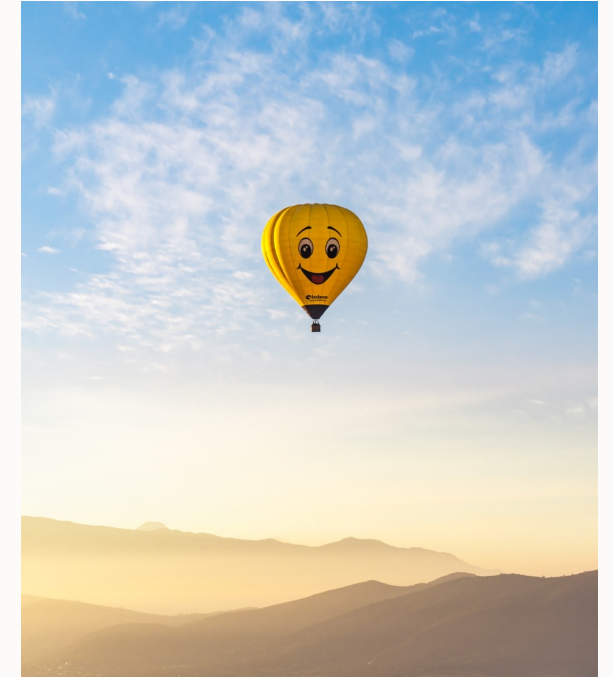
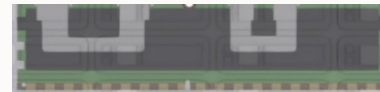
- Rebalance happens due to disk failure
- Primary mirror was cached in Data Accelerator
- Primary mirror goes to other cell
- Cache in Data Accelerator follows
- Latency preserved end-user happy



Data
Accelerator



Data
Accelerator



Credit Jacob Vizek
<https://unsplash.com/photos/ibvHQnpk4LE>



Exadata : Brownout

Cisco RoCE Spine Switch Software Update

- MAA team tested spine switch reboots in Multi Rack configurations
- Exadata 21.2.*
 - No blackouts
 - Significantly reduced brownout



Lifecycle Management

Data Protection



Brownout

Quality Of Service and Performance



Exadata : Life Cycle Management

Exachk

- Recommendations from Exachk come straight from the field & engineering discussed in weekly meetings
- It is crucial to always have the latest version
 - Keeps track of critical issues
 - Issues in certain releases
 - We don't allow to run if older than 180 days
- Highly recommended to run Exachk
 - Once a month
 - Before and after any major configuration change eg : patching, storage addition
- Best practice health check

Database Server		
Status	Type	Message
CRITICAL	Database Check	Database parameter CLUSTER_INTERCONNECTS is not set to the recommended value
CRITICAL	Database Check	Database parameters log_archive_dest_n with Location attribute are not all set to recommended value
CRITICAL	OS Check	Hardware and firmware profile check is not successful. [Database Server]
CRITICAL	OS Check	The InfiniBand Address Resolution Protocol (ARP) Configuration on Database Servers should be as recommended
FAIL	SQL Check	Some data or temp files are not autoextensible
FAIL	OS Check	Memlock settings do not meet the Oracle best practice recommendations
FAIL	ASM Check	Fast recovery area allocation totals are greater than the total space of the DB_RECOVERY_FILE_DEST disk group
FAIL	OS Check	Active kernel version should match expected version for installed Exadata Image
FAIL	OS Check	One or more database server has non-test stateless alerts with null "examinedby" fields
FAIL	OS Check	One or more database servers have stateful alerts that have not been cleared
FAIL	Database Check	Hidden database Initialization Parameter usage is not correct
WARNING	Database Check	Local listener init parameter is not set to local node VIP
WARNING	Database Check	Database parameter DB_BLOCK_CHECKING on PRIMARY is NOT set to the recommended value.
INFO	OS Check	Exadata Critical Issues (Doc ID 1270094.1):- DB1-DB4,DB6,DB9-DB41, EX1-EX54,EX56 and IB1-IB3,IB5-IB8
INFO	Database Check	One or more non-default AWR baselines should be created



Exadata : Life Cycle Management

Exachk

Cluster Summary

Cluster Name	Cluster-c1
OS/Kernel Version	LINUX X86-64 OELRHHEL 7 4.14.35-2047.505.4.4.el7uek.x86_64
CRS Home - Version	/u01/app/21.0.0.0/grid - 21.3.0.0.0
DB Home - Version - Names	/u01/app/oracle/product/21.0.0.0/dbhome_1 - 21.3.0.0.0 - cdm213 database /u01/app/oracle/product/19.0.0.0/dbhome_1 - 19.12.0.0.0 - cdm19c database /u01/app/oracle/product/18.0.0.0/dbhome_1 - 18.14.0.0.0 - cdm18c database /u01/app/oracle/product/12.2.0.1/dbhome_1 - 12.2.0.1.210710 - cdm122 database /u01/app/oracle/product/12.1.0.2/dbhome_1 - 12.1.0.2.210710 - 3 databases
Exadata Version	21.2.4.0.0
Number of nodes	8
Database Servers	2
Storage Servers	3
IB Switches	3
EXAchk Version	21.3.0_20211029
Collection	exachk_random01client01_rac12c_030922_00257
Duration	32 mins, 6 seconds
Executed by	root
Arguments	-hardwaretype X4-2
Collection Date	10-Mar-2022 00:57:54

- There are 6 flagged critical checks, 16 flagged failed checks , 7 flagged warning checks, 17 flagged info checks. By default, 16 failed checks are considered critical.
- This version of EXAchk is considered valid for 48 days from today or until a new version is available

Exadata Critical Issues

The following Exadata Critical Issues ([MOS Note 1270094.1](#)) have been checked in this report:

- Exadata Storage Server : EX1-EX65,EX67,EX69,EX70
- Database Server : DB1-DB4, DB6, DB9-DB49
- InfiniBand switch : IB1-IB3,IB5-IB9

Cluster Summary

Cluster Name	Cluster-c1
OS/Kernel Version	LINUX X86-64 OELRHHEL 7 4.14.35-2047.505.4.4.el7uek.x86_64
CRS Home - Version	/u01/app/21.0.0.0/grid - 21.3.0.0.0
DB Home - Version - Names	/u01/app/oracle/product/21.0.0.0/dbhome_1 - 21.3.0.0.0 - cdm213 database /u01/app/oracle/product/19.0.0.0/dbhome_1 - 19.12.0.0.0 - cdm19c database /u01/app/oracle/product/18.0.0.0/dbhome_1 - 18.14.0.0.0 - cdm18c database /u01/app/oracle/product/12.2.0.1/dbhome_1 - 12.2.0.1.210720 - cdm122 database /u01/app/oracle/product/12.1.0.2/dbhome_1 - 12.1.0.2.210720 - 3 databases
Exadata Version	21.2.4.0.0
Number of nodes	8
Database Servers	2
Storage Servers	3
IB Switches	3
EXAchk Version	21.4.2_20220211
Collection	exachk_random01client01_rac12c_030922_151642
Duration	30 mins, 48 seconds
Executed by	root
Arguments	-hardwaretype X4-2
Collection Date	09-Mar-2022 15:22:38

- There are 5 flagged critical checks, 19 flagged failed checks , 6 flagged warning checks, 18 flagged info checks. By default, 19 failed checks are considered critical.
- This version of EXAchk is considered valid for 154 days from today or until a new version is available

Exadata Critical Issues

The following Exadata Critical Issues ([MOS Note 1270094.1](#)) have been checked in this report:

- Exadata Storage Server : EX1-EX65,EX67,EX69,EX70,EX71,EX72
- Database Server : DB1-DB4, DB6, DB9-DB49
- InfiniBand switch : IB1-IB3,IB5-IB9



Exadata : Life Cycle Management

Exachk: top observed painpoints

- Huge pages not set correctly
 - In later DB releases we check if SGA > 32 Gb is used without huge pages configured
 - If that is the case the instance doesn't start
- Redundancy recommendation not followed
- Critical issue that is already fixed in later releases

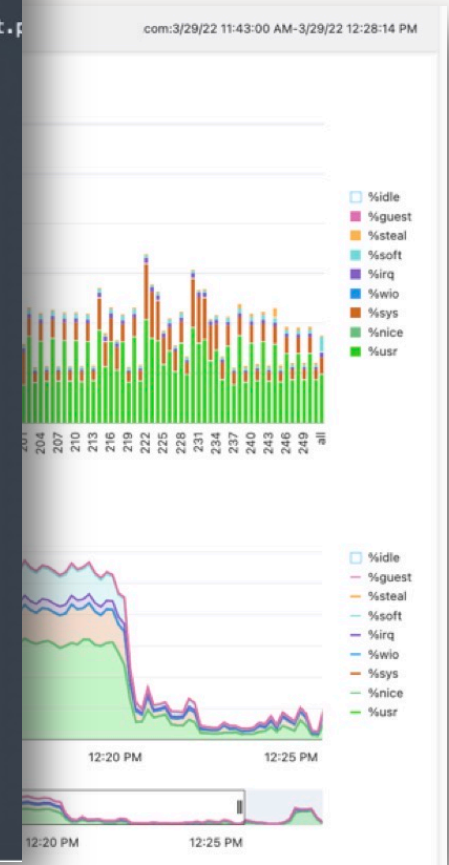


Exadata : Life Cycle Management

Exawatcher : Graphing

```
zzz <08/26/2021 01:07:37> subcount: 1
top - 01:07:39 up 23 days, 11:34, 0 users, load average: 7.01, 7.94, 8.74
Threads: 14011 total, 5 running, 11133 sleeping, 0 stopped, 0 zombie
%Cpu(s): 4.8 us, 2.1 sy, 0.0 ni, 92.6 id, 0.0 wa, 0.2 hi, 0.3 si, 0.0 st
KiB Mem : 15834616+total, 22833467+free, 10630640+used, 29206297+buff/cache
KiB Swap: 25165820 total, 25158640 free, 7180 used. 45395289+avail Mem

  PID USER      PR  NI   VIRT   RES    SHR  S %CPU  %MEM    TIME+  COMMAND
 31661 splunk   20   0  185704 12144  7236  R 99.9  0.0   0:07.23 python /opt/splunkforwarder/etc/apps/bk_elp_update_read_perm/bin/generate_monitored_list.p
35530 exawatch 20   0  177024 19592  4188  R 80.0  0.0   0:00.72 /usr/bin/top -b -c -H -n 1 -w512
35561 exawatch 20   0  169680 12436  4160  R 40.0  0.0   0:00.32 /usr/bin/top -b -c -n 1 -w512
    1 root     20   0  217384 11300  6244  R 20.0  0.0   7669:00 /usr/lib/systemd/systemd --switched-root --system --deserialize 22
257794 oracle  20   0   19.0g 156488 139184 S  7.5  0.0   6:13.84 oracle (LOCAL=NO)
137634 grid   20   0  3658556 104120  98748 S  5.0  0.0 156:28.75 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
229236 grid   20   0  3658556 110608 105236 S  5.0  0.0 156:45.62 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
 35565 oracle  20   0  7050676 69780  65252 S  3.8  0.0   0:00.03 ora_
110869 grid   20   0  3654816 80352  76448 S  3.8  0.0 33:29.37 asm_pmon+ASM4
137609 grid   20   0  3658560 113088 107716 S  3.8  0.0 157:26.15 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
137647 grid   20   0  3659408 112392 106304 S  3.8  0.0 158:49.41 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
172236 grid   20   0  3658556 107148 101720 S  3.8  0.0 10:25.29 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
172284 grid   20   0  3658560 103324  97956 S  3.8  0.0 10:20.94 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
178226 grid   20   0  3658564 104116  98748 S  3.8  0.0 10:20.00 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
192223 grid   20   0  3668324 113380 107376 S  3.8  0.0 14:48.76 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
225979 grid   20   0  3658560 118500 113124 S  3.8  0.0 163:44.29 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
226001 grid   20   0  3658556 105780 100404 S  3.8  0.0 163:30.96 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
226633 grid   20   0  3658560 111036 105680 S  3.8  0.0 164:25.67 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
227118 grid   20   0  3659412 113076 106964 S  3.8  0.0 161:17.46 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
227320 grid   20   0  3658564 114312 108940 S  3.8  0.0 162:27.96 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
227628 grid   20   0  3658560 118404 113028 S  3.8  0.0 162:31.65 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
228001 grid   20   0  3658560 112884 107512 S  3.8  0.0 164:01.06 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
228157 grid   20   0  3658564 106080 100704 S  3.8  0.0 163:24.08 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
229276 grid   20   0  3658564 103220  97840 S  3.8  0.0 158:57.30 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
230200 grid   20   0  3658564 117016 111640 S  3.8  0.0 165:29.22 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
230209 grid   20   0  3658564 116896 111520 S  3.8  0.0 165:27.66 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
249654 grid   20   0  3668324 110712 104708 S  3.8  0.0 57:58.89 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
253216 grid   20   0  3669152 109952 102936 S  3.8  0.0 10:11.55 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
254595 grid   20   0  3669160 114024 107272 S  3.8  0.0 159:07.87 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
255128 grid   20   0  3658556 109576 104200 S  3.8  0.0 159:17.33 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
257510 grid   20   0  3658556 104528  99160 S  3.8  0.0 10:06.37 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
257528 grid   20   0  3659388 105640  99552 S  3.8  0.0 10:01.25 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
257568 grid   20   0  3659408 104836  98740 S  3.8  0.0 9:55.44 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
264924 grid   20   0  3669156 120184 113128 S  3.8  0.0 156:47.16 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
307188 grid   20   0  3668328 109736 103732 S  3.8  0.0 12:38.59 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
  2185 grid   20   0  3668332 108372 102356 S  2.5  0.0 12:04.51 oracle+ASM4_ (DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq)))
```



Exadata : Life Cycle Management

Monitoring by Enterprise Manager 13c

ORACLE Enterprise Manager Cloud Control 13c

Exadata Storage Server

Status: ↑ Up
Health: ✓ Good
IORM Status: Enabled
Release Version: 21.2.9.0.0

Capacity

Disk Size: 0.00%
Flash (GB): 47,693 (95%)

ASM Disk Group Capacity

DATA1	36232 GB
RECO1	9075 GB

Performance

IO Queue Per Disk

Average Flash Response Time

I/O Activity

Drive Type	Utilization (%)	IO/Sec	MB/Sec
Flash Drive	0	69	93.9

Flash Cache: Read Hit (%) 0, Write IOPS 0

Smart Scan IO (MB/Sec): Flash 0.4, Storage Index 0

Infiniband (MB/Sec): Network Sent 0, Network Received 0

Smart Log Efficiency (%) 100, IORM Boost 0

I/O Distribution by Databases

I/O Utilization (%)

Database	Utilization (%)
MAA19	~100
ASM	~100
ABTSTDB	~100
CORPDB19	~100
WF198CDB	~100
Others	~100

Flash I/O Service Time (ms/request)

Database	Flash I/O Time (ms)	Flash IORM Wait Time (ms)
MAA19	~0.12	~0.02
ASM	~0.10	~0.02
ABTSTDB	~0.10	~0.02
CORPDB19	~0.10	~0.02
WF198CDB	~0.10	~0.02
Others	~0.15	~0.02

Incidents and Problems

No matching incidents or problems found.

Enterprise

View: Schematic, Photo-Realistic, Table

Zoom: 100%, 50%

Show: Temperature, Empty Slot

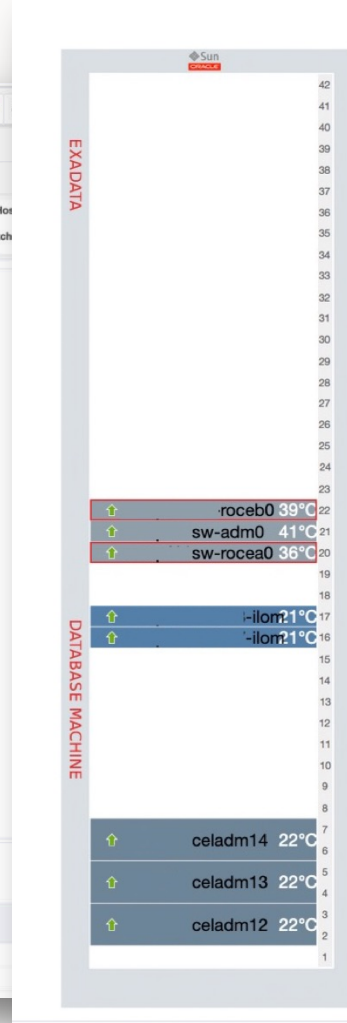
Status: Up, Down, Blackout, Locator Light

Component: Ethernet Switch, Server, Storage Server

Temperature Readings:

- roceb0: 39°C
- sw-adm0: 41°C
- sw-rocea0: 36°C
- ilom2: 21°C
- ilom1: 21°C
- celadm14: 22°C
- celadm13: 22°C
- celadm12: 22°C

Target	Severity	Status
Target 1	High	New
Target 2	High	New



View

- Schematic
- Photo-Realistic
- Table

Zoom

- 100%
- 50%

Show

- Temperature
- Empty Slot

Status

- ↑ Up
- ↓ Down
- ✖ Blackout
- 📶 Locator Light

Component

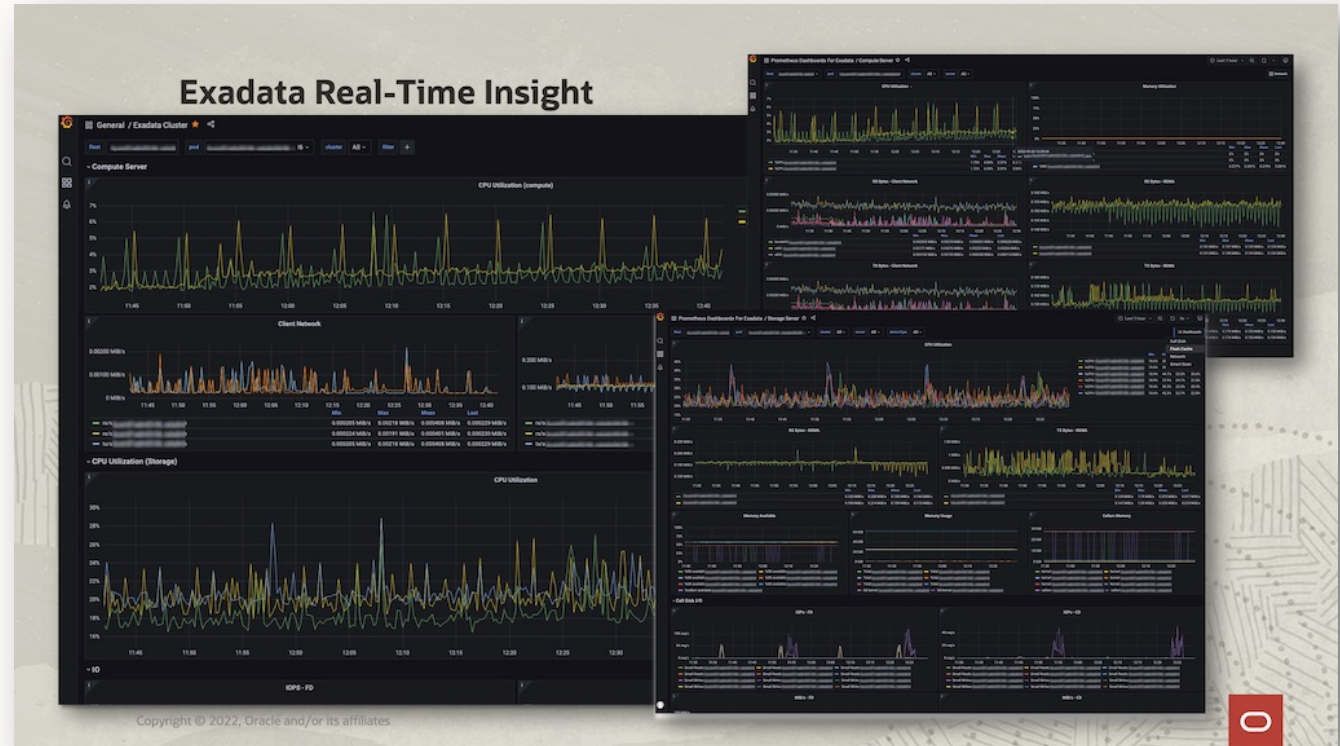
- Ethernet Switch
- Server
- Storage Server



Exadata : Life Cycle Management

Exadata Real Time Insight

- Automatically stream up-to-the-second metric observations from all servers in your Exadata fleet
- Feed customizable monitoring dashboards for real-time analysis and problem-solving
- Comprehensive : > 200 Exadata Soft- & Hardware Metrics
- Proactive issue detection and real-time decision making



Exadata : Life Cycle Management

Exadata Real-Time Insight

- Automatically stream up-to-the-second metric observations from all servers in your Exadata fleet
- Feed customizable monitoring dashboards for real-time analysis and problem-solving

• Comprehensive

- 200+ Exadata Software & Hardware Metrics
- Fine-grained metrics can be collected as often as every 1 second

• Integrated

- Integrated with popular time-series and observability platforms
- Stream fine-grained metrics to user-defined endpoints in real time

• Insightful

- Enables proactive issue detection and real-time decision making



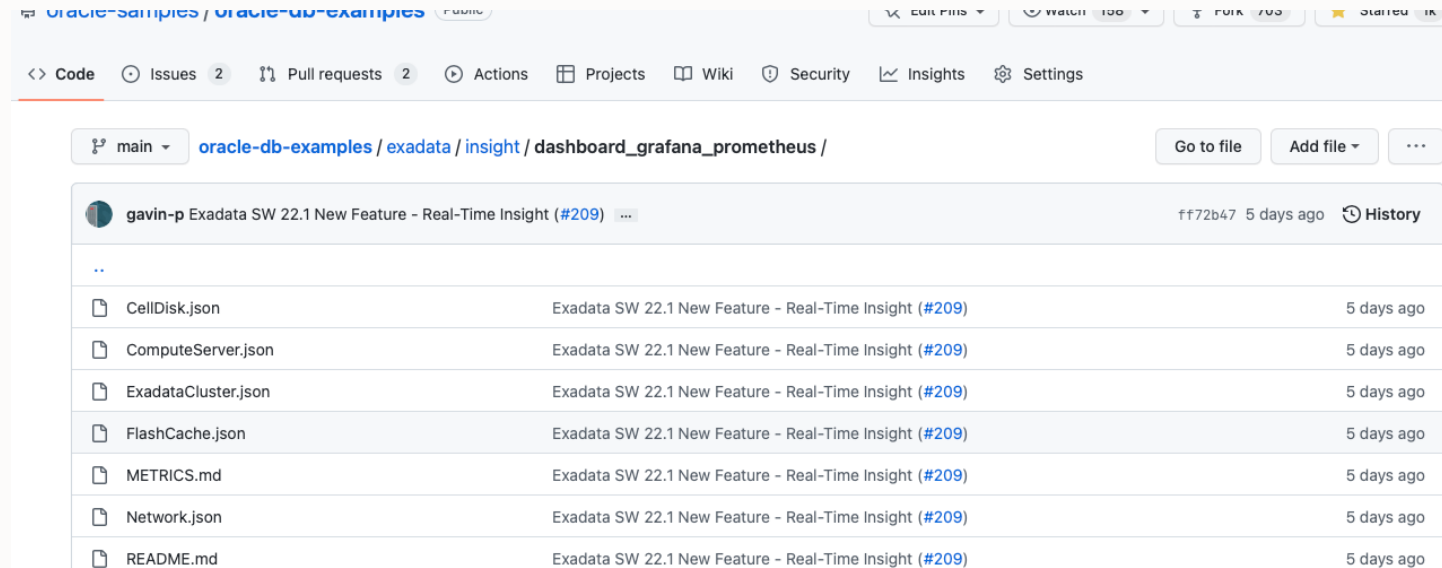
<https://blogs.oracle.com/exadata/post/real-time-insight-quick-start>



Exadata : Life Cycle Management

Exadata Real-Time Insight – Sample Dashboards Code

- Oracle Samples repository on GitHub.com contains example Real-Time Insight dashboards.
 - The following dashboard code is included (Grafana/Prometheus):
 - Exadata Cluster
 - Compute
 - Storage Server
 - Cell Disk
 - Flash Cache
 - Smart Scan
 - Network



- <https://github.com/oracle-samples/oracle-db-examples/tree/main/exadata/insight>



Exadata : Life Cycle Management

Exadata Real-Time Insight – Sample Dashboards

- **Exadata Cluster:** Provides a cluster-wide view that shows metrics for compute nodes and storage servers
- **Compute:** Provides a compute-node view that shows CPU and network utilization for the compute nodes
- **Storage Server:** Provides a storage-server-centric view that focuses on storage server CPU and I/O metrics, as well as Exadata metrics for Smart Flash Cache, Smart Flash Log, and Smart I/O
 - **Cell Disk:** Shows cell disk I/O metrics on the storage server
 - **Flash Cache:** Shows flash cache metrics on the storage server
 - **Smart Scan:** Shows smart scan metrics on the storage server



Exadata : Life Cycle Management

Exadata AWR support

Exadata Configuration and Statistics

- [Exadata Report Summary](#)
- [Exadata Server Configuration](#)
- [Exadata Server Health Report](#)
- [Exadata Statistics](#)

[Back to Top](#)

Exadata Server Configuration

- [Exadata Storage Server Model](#)
- [Exadata Storage Server Version](#)
- [Exadata Storage Information](#)
- [Exadata Griddisks](#)
- [Exadata Celldisks](#)
- [ASM Diskgroups](#)
- [IORM Objective](#)

Exadata Storage Server Model

- Model Information of Servers
- CPU Count refers to logical CPUs, including cores and hyperthreads

Model	CPU Count	Memory (GB)	# Cells
Oracle Corporation ORACLE SERVER X9-2L High Capacity	96/96	252	12

[Back to Exadata Server Configuration](#)

Exadata Storage Server Version

- Version information of packages on the storage server

Package Type	Package Version	Cells
Kernel	4.14.35-2047.511.5.5.el7uek.x86_64	All (12)
Cell	cell-21.2.10.0.0_LINUX.X64_220317-1.x86_64	All (12)
Offload	celloff-11.2.3.3.1_LINUX.X64_200526	All (12)
Offload	celloff-21.2.10.0.0_LINUX.X64_220317	All (12)
Offload	celloff-12.1.2.4.0_LINUX.X64_210708.1	All (12)

[Back to Exadata Server Configuration](#)

Exadata Storage Information

- Storage information per cell
- 'Total' is the sum for all cells

# Cells	Size (GB)			# Celldisks			# Griddisks	Cell Name
	Flash Cache	PMEM Cache	Flash Log	Hard Disk	Flash	PMEM		
11	23,845.81	1,500.56	0.50	12	4	12	24 (11):	
1	23,845.81	1,500.56	0.38	12	4	12	24 (1):	
Total (12)	286,149.75	18,006.75	5.88	144	48	144	288 All (12)	

[Back to Exadata Server Configuration](#)

Exadata Griddisks

- Griddisks on the storage servers
- Disk Type <F|H|M><size>: F-Flash, H-Harddisk, M-persistent Memory
- Size (GB) - Griddisk: indicates size of individual Griddisks in the cells
- Size (GB) - Cell Total: indicates total size per cell
- Size (GB) - System Total: indicates total size over all cells

Outlier Summary - Disk Level

- Outliers are disks whose average performance is outside the normal range, where normal range is +/- 3 standard deviation
- Outlier disks must have a minimum of 10 IOPs. Idle disks are not considered for outlier analysis.
- Outliers for hard disks are displayed when Hard Disk IOPs exceeds 306.72 (1% of maximum capacity of 30,672)
- Outliers for flash disks are displayed when Flash Disk IOPs exceeds 95260.8 (1% of maximum capacity of 9,526,080)
- Outliers for flash disks will not be displayed. There are only 86,566 flash IOPs
- A 'v' and a dark yellow background indicates an outlier value below the low range
- A '^' and a light red background indicates an outlier value above the high range
- A '*' and a dark red background indicates over maximum capacity
- Disk Type <F|H|M><size>: F-Flash, H-Hard Disk, M-pMEM; PMEM I/O only include remote I/Os processed by cellsvr
- Maximum hard disk capacity: IOPS: H/16.0T: 213 | IO MB/s: H/16.0T: 148
- Maximum flash disk capacity: IOPS: F/5.8T: 198,460 | IO MB/s: F/5.8T: 11,250

Disk Name	Cell Name	Disk Type	Statistic Name	Value	Mean	Std Dev	Normal Range
CD_01_		H/16.0T	% Disk Utilization	^ 2.83	0.34	0.61	0.00 - 2.16
CD_02_		H/16.0T	Small Reads/s	^ 39.95	3.37	7.43	0.00 - 25.67
CD_06_		H/16.0T	Small Reads/s	^ 39.44	3.37	7.43	0.00 - 25.67
CD_08_		H/16.0T	Small Reads/s	^ 39.24	3.37	7.43	0.00 - 25.67
CD_08_		H/16.0T	Small Reads/s	^ 39.39	3.37	7.43	0.00 - 25.67
CD_10_		H/16.0T	Small Reads/s	^ 39.83	3.37	7.43	0.00 - 25.67
CD_10_		H/16.0T	% Disk Utilization	^ 2.45	0.34	0.61	0.00 - 2.16

[Back to Exadata Outlier Summary](#)

[Back to Exadata Resource Statistics](#)



Exadata : Life Cycle Management

Custom Diagnostic Package for Storage Server Alerts

[.oracle.com](#)

[.]_diag

- [alert.log \(Full version\)](#)
- [ms-odl-316.trc](#)
- [ms-odl-317.trc](#)
- [ms-odl-318.trc](#)
- [ms-odl-319.trc](#)
- [ms-odl-320.trc](#)
- [ms-odl-321.trc](#)
- [ms-odl-322.trc](#)
- [ms-odl-323.trc](#)
- [ms-odl-324.trc](#)
- [ms-odl-325.trc](#)
- [ms-odl-326.trc](#)
- [ms-odl-327.trc](#)
- [ms-odl-328.trc](#)
- [ms-odl-329.trc](#)
- [ms-odl-330.trc](#)
- [ms-odl-331.trc](#)
- [ms-odl-332.trc](#)
- [ms-odl-333.trc](#)
- [ms-odl-334.trc \(Full version\)](#)
- [ms-odl.trc \(Full version\)](#)

[.]_var

- [\[+\] log](#)

[.]_ExaWatcher

```

/cell/cellsrv/deploy/config/metadata/5e901e50-6dc7-4d34-9928-4af32307d502)
2022-05-11T09:25:22.228533-07:00
System Disk metadata update info: DATAC1_CD_01_ : celldisk update for cachedby list succeeded
2022-05-11T09:25:22.239524-07:00
System Disk metadata update info: DATAC1_CD_07_ : celldisk update for cachedby list succeeded
2022-05-11T09:25:22.286516-07:00
System Disk metadata update info: DATAC1_CD_05_ : celldisk update for cachedby list succeeded
2022-05-11T09:25:22.443236-07:00
System Disk metadata update info: DATAC2_CD_01_ : celldisk update for cachedby list succeeded
2022-05-11T09:25:22.639985-07:00
System Disk metadata update info: DATAC1_CD_03_ : celldisk update for cachedby list succeeded
2022-05-11T09:25:22.695806-07:00
System Disk metadata update info: RECOC1_CD_01_ : celldisk update for cachedby list succeeded
2022-05-11T09:25:22.776889-07:00
System Disk metadata update info: DATAC2_CD_07_ : celldisk update for cachedby list succeeded
2022-05-11T09:25:22.850375-07:00
System Disk metadata update info: DATAC2_CD_05_ : celldisk update for cachedby list succeeded
2022-05-11T09:25:23.053280-07:00
System Disk metadata update info: RECOC1_CD_07_ : celldisk update for cachedby list succeeded
2022-05-11T09:25:23.122820-07:00
System Disk metadata update info: RECOC1_CD_05_ : celldisk update for cachedby list succeeded
2022-05-11T09:25:23.228778-07:00
System Disk metadata update info: DATAC2_CD_03_ : celldisk update for cachedby list succeeded
2022-05-11T09:25:23.270691-07:00
System Disk metadata update info: RECOC2_CD_01_ : celldisk update for cachedby list succeeded
2022-05-11T09:25:23.284951-07:00
System Disk metadata update info: RECOC2_CD_07_ : celldisk update for cachedby list succeeded
2022-05-11T09:25:23.762756-07:00
System Disk metadata update info: RECOC2_CD_05_ : celldisk update for cachedby list succeeded
2022-05-11T09:25:24.415166-07:00
System Disk metadata update info: RECOC2_CD_03_ : celldisk update for cachedby list succeeded
[MS] Disk controller was hung. Cell was power cycled to restore access to the cell. Timestamp: Wed May 11 09:11:17 PDT 2022
```

Maintenance: Hardware Alert 64

Event Time 2022-05-11T09:11:17-07:00

Description Disk controller was hung. Cell was power cycled to restore access to the cell.

Affected Cell

Name	
Server Model	Oracle Corporation ORACLE SERVER X8-2L
Chassis Serial Number	
Release Version	22.1.0.0.0.220504
RPM Version	22.1.0.0.0_LINUX.X64_220504-1

Recommended Action Informational.

succeeded

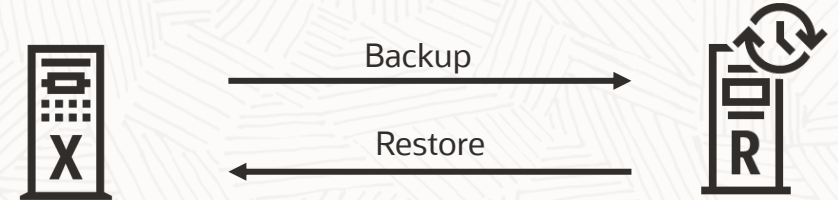
May 11 09:11:17 PDT 2022



Exadata : Life Cycle Management

Backup

- Backup your databases 😊
 - ZDLRA or ZFS appliance with RMAN are recommended
- Backup KVM host and KVM Guest
 - More details see at end of presentation
- Test your backups



Schrödinger Backup :

The condition of any backup is unknown until a restore is attempted



Exadata Live Update

Increase security and minimize database server and VM reboots

Exadata System Software provides operating system, firmware, and Exadata software updates that are crucial for the optimal and secure operation of Exadata Database Servers and Oracle Database

Updates are applied in a rolling fashion across database servers

Exadata Live Update applies updates online and defers any remaining work to occur at a scheduled time

Exadata Live Update uses familiar Linux technologies, including RPM and ksplice, to **apply updates online** to database servers/VMs **avoiding the need to reboot**



Exadata Live Update Options

Exadata Live Update multiple options based on the Common Vulnerability Scoring System (CVSS). When using Exadata Live Update, you choose from the following options:

- highcvss** Applies only security updates to address vulnerabilities with a CVSS score of 7 or greater
- allcvss** Applies only security updates to address vulnerabilities with any CVSS score
- full** Performs a full update, which includes all security-related updates and all other non-security updates. Equivalent to regular updates applied with a server/VM reboot

```
$ patchmgr --dbnodes kvm_guests.lst --upgrade --repo <repo.zip location> --rolling \  
--target_version 24.1.0.0.0.240517 --live-update-target highcvss|allcvss|full
```



Viewing Outstanding Work

Not all update content can be applied online, or activated without a reboot

- e.g. firmware, booting with the latest kernel, JDK

These updates are called ‘outstanding work’ and are staged for activation at the next **graceful** shutdown

Use `patchmgr --live-update-list-outstanding-work` to show outstanding items

```
$ patchmgr --dbnodes kvm_guests.lst --live-update-list-outstanding-work
***
Summary of outstanding work for Exadata Live Update:
exdpm1adm01vm01.example.com: (*) 2024-08-15 00:17:08: Exadata Live Update outstanding work is
scheduled for completion at the next reboot
    - The Linux kernel will be updated from version 5.4.17-
      2136.330.7.5.el8uek to 5.4.17-2136.333.5.1.el8uek.
      Current Uptrack kernel version: 5.4.17-2136.333.5.1.el8uek.x86_64
    - New package uptrack-updates-5.4.17-2136.333.5.1.el8uek.x86_64
      (version 20240725-0) will be installed.
```

Applying Outstanding Work

By default, outstanding work is applied during the next graceful shutdown

Administrators can use `patchmgr --live-update-schedule-outstanding-work` to

- Specify the reboot window - "YYYY-MM-DD HH24:MM:SS TZ"

```
$ patchmgr --dbnodes kvm_guests.lst --live-update-schedule-outstanding-work \  
          "2024-11-04 22:00:00 AEDT"
```

- To defer applying outstanding work – ‘never’

```
$ patchmgr --dbnodes kvm_guests.lst --live-update-schedule-outstanding-work never
```

- Reset a previously set schedule to the default behavior

```
$ patchmgr --dbnodes kvm_guests.lst --live-update-schedule-outstanding-work reset
```

Oracle recommends outstanding work be applied at least every 3 months



Exadata Live Update Best Practices

Database Server/VM Backup

- Patchmgr automatically creates a system backup during all updates to allow for fast rollback if required
- Additional administrator-managed backups are recommended

Graceful reboots

- Include `vm_maker --stop_domain/--start_domain` operations, host restart (`shutdown -r`), a short press of the power button on the server, etc.
- Restarting the physical database server also restarts VMs
 - Useful (but not required) to align VM and physical server reboot
- Avoid resetting VMs and physical servers while outstanding work is applied

Use Database MAA features including Transparent Application Continuity to mask planned reboot from applications and users

Exadata Live Update

Applying monthly maintenance releases - examples

Quarterly Update Windows (Recommended)

August	September	October	November
<ul style="list-style-type: none"> • 24.1.3 • Full Update • Server/VM reboot 	<ul style="list-style-type: none"> • 24.1.4 • Exadata Live Update • No reboot 	<ul style="list-style-type: none"> • 24.1.5 • Exadata Live Update • No reboot 	<ul style="list-style-type: none"> • 24.1.6 • Full Update • Server/VM reboot
December	January	February	March
<ul style="list-style-type: none"> • 24.1.7 • Exadata Live Update • No reboot 	<ul style="list-style-type: none"> • 24.1.8 • Exadata Live Update • No reboot 	<ul style="list-style-type: none"> • 24.1.9 • Full Update • Server/VM reboot 	<ul style="list-style-type: none"> • 24.1.10 • Exadata Live Update • No reboot

Bi-Yearly Update Windows

August	September	October	November
<ul style="list-style-type: none"> • 24.1.3 • Full Update • Server/VM reboot 	<ul style="list-style-type: none"> • 24.1.4 • Exadata Live Update • No reboot 	<ul style="list-style-type: none"> • 24.1.5 • Exadata Live Update • No reboot 	<ul style="list-style-type: none"> • 24.1.6 • Exadata Live Update • No reboot
December	January	February	March
<ul style="list-style-type: none"> • 24.1.7 • Exadata Live Update • No reboot 	<ul style="list-style-type: none"> • 24.1.8 • Exadata Live Update • No reboot 	<ul style="list-style-type: none"> • 24.1.9 • Full Update • Server/VM reboot 	<ul style="list-style-type: none"> • 24.1.10 • Exadata Live Update • No reboot

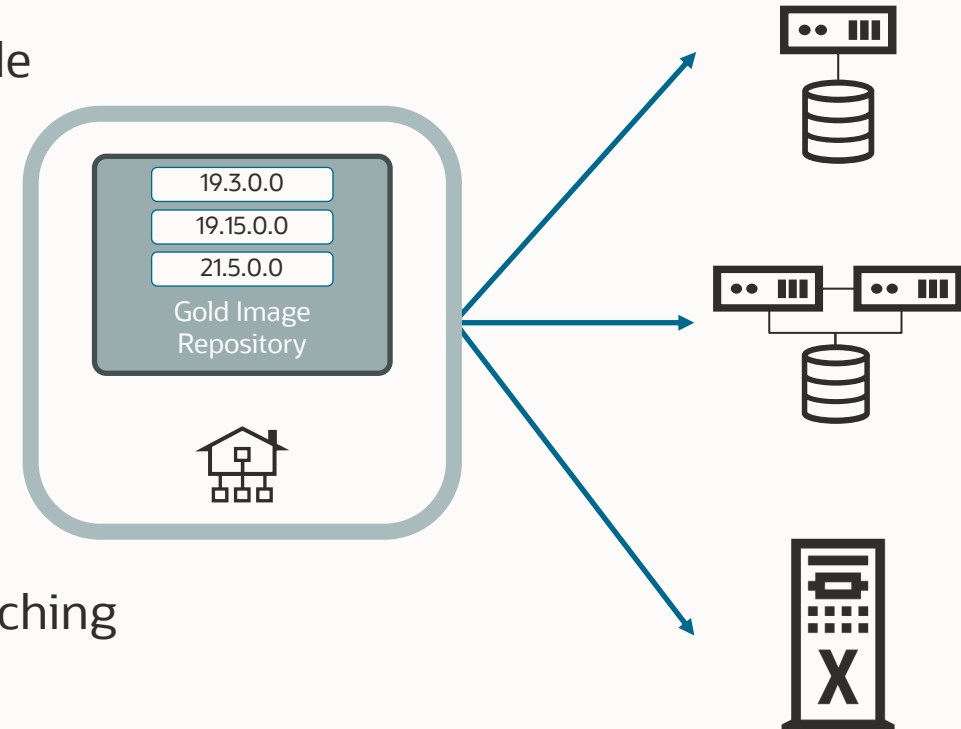


Exadata : Life Cycle Management

Planned Maintenance

Exadata patchmgr utility can be used to patch the whole hardware stack :

- Storage cells
- RoCE switches
- Admin switches
- Baremetal and KVM Host
- KVM Guest



Fleet Patching & Provisioning the tool for out place patching

- Database homes
- Grid infrastructure and combined GI + DB patching
- Also Exadata patching
- www.oracle.com/goto/FPP
- One tool to patch / upgrade your whole Oracle DB stack



Exadata : Further Reading

Backup

- <https://www.oracle.com/technetwork/database/availability/recovery-appliance-maint-practices-4487388.pdf>

KVM Virtualization

- <https://www.oracle.com/a/tech/docs/exadata-kvm-overview.pdf>

Life Cycle Management

- <https://www.oracle.com/a/tech/docs/exadata-software-maintenance-2022.pdf>

Security

- <https://www.oracle.com/a/tech/docs/exadata-maximum-security-architecture.pdf>

Exadata Real-Time Insight

- <https://blogs.oracle.com/exadata/post/exadata-real-time-insight>



Reference

Useful Resources

Exadata Product Management Blog - <https://blogs.oracle.com/exadata/>

MOS Note Reference Blog - <https://blogs.oracle.com/exadata/post/exadata-mos-notes>

Exadata Database Machine and Exadata Storage Server Supported Versions (Doc ID [888828.1](#))

Oracle Exadata Database Machine EXAchk (Doc ID [1070954.1](#))

Oracle Exadata Best Practices (Doc ID [757552.1](#))

Exadata Critical Issues (Doc ID [1270094.1](#))

Exadata Patching Overview and Patch Testing Guidelines (Doc ID [1262380.1](#))

The ASM Priority Rebalance feature - An Example (Doc ID [1968607.1](#))

Physical and Logical Block Corruptions. All you wanted to know about it. (Doc ID [840978.1](#))

Best Practices for Corruption Detection, Prevention, and Automatic Repair - in a Data Guard Configuration (Doc ID [1302539.1](#))

Understanding ASM Capacity and Reservation of Free Space in Exadata (Doc ID [1551288.1](#))



Exadata MAA : Conclusion

Solid as a rock



Credit : Zoltan Tasi <https://unsplash.com/photos/QxjEi8Fs9Hg>

Out of this world performance



Credit : Space X <https://unsplash.com/photos/OHOU-5UVIYQ>



Thank you



Our mission is to help people see
data in new ways, discover insights,
unlock endless possibilities.



ORACLE