# "Extending IN-Memory Database Processing to Shared Flash

**Gurmeet Goindi**
Master Product Manager

# Safe Harbor Statement

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.

# Exadata Database Machine

Performance, Availability and Security

**Best Platform for Oracle Databases on-premises and in the Cloud**
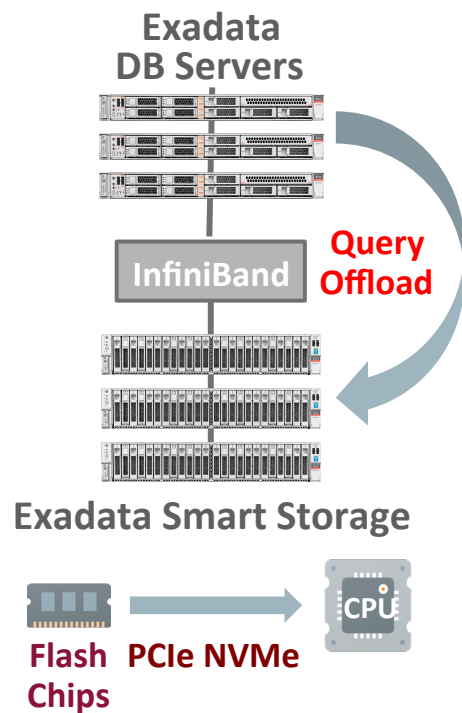
Enabled by:

- Single-vendor accountability

- Exclusive focus on databases

- Deep h/w and s/w integration

- Revolutionary approach to storage

ORACLE®

# Exadata Achieves Memory Performance with Shared Flash

**Exadata DB Servers**



**InfiniBand**

**Query Offload**

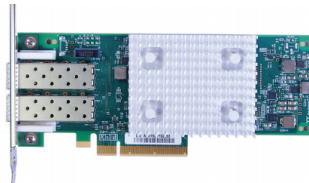**Exadata Smart Storage**

**Flash Chips** **PCIe NVMe** → CPU

- Exadata X6 delivers 300GB/sec flash bandwidth to <u>any</u> server
  - Approaches 800GB/sec aggregate DRAM bandwidth of DB servers
- Must move compute to data to achieve full flash potential
  - Requires owning full stack, can't be solved in storage alone
- Fundamentally, storage arrays can share flash <u>capacity</u> but not flash <u>performance</u>
  - Even with next gen scale-out, PCIe networks, or NVMe over fabric
  - e.g. new EMC DSSD has 3-6 times lower throughput than Exadata X6
- Shared storage with memory-level bandwidth is a paradigm change in the industry
  - Get near DRAM throughput, with the capacity of shared flash

# NVMe PCI-e Flash Disrupts the Storage Array Model
## New improvements are causing 100X bottlenecks across shared storage stack
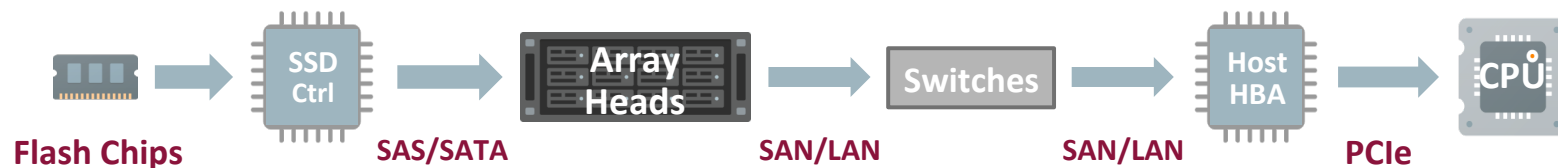
**Latest PCIe Flash**
5.4 GB/sec

**SAN Link = 40Gb**
5 GB/sec
Less than 1 Flash card

**Leading All Flash Array**
24 GB/sec
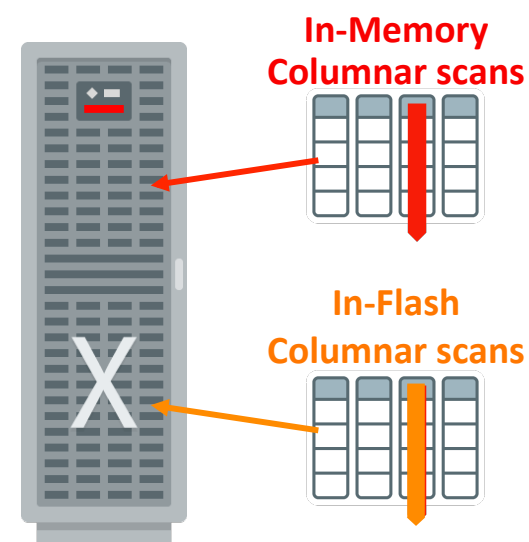Less than 5 Flash card

---

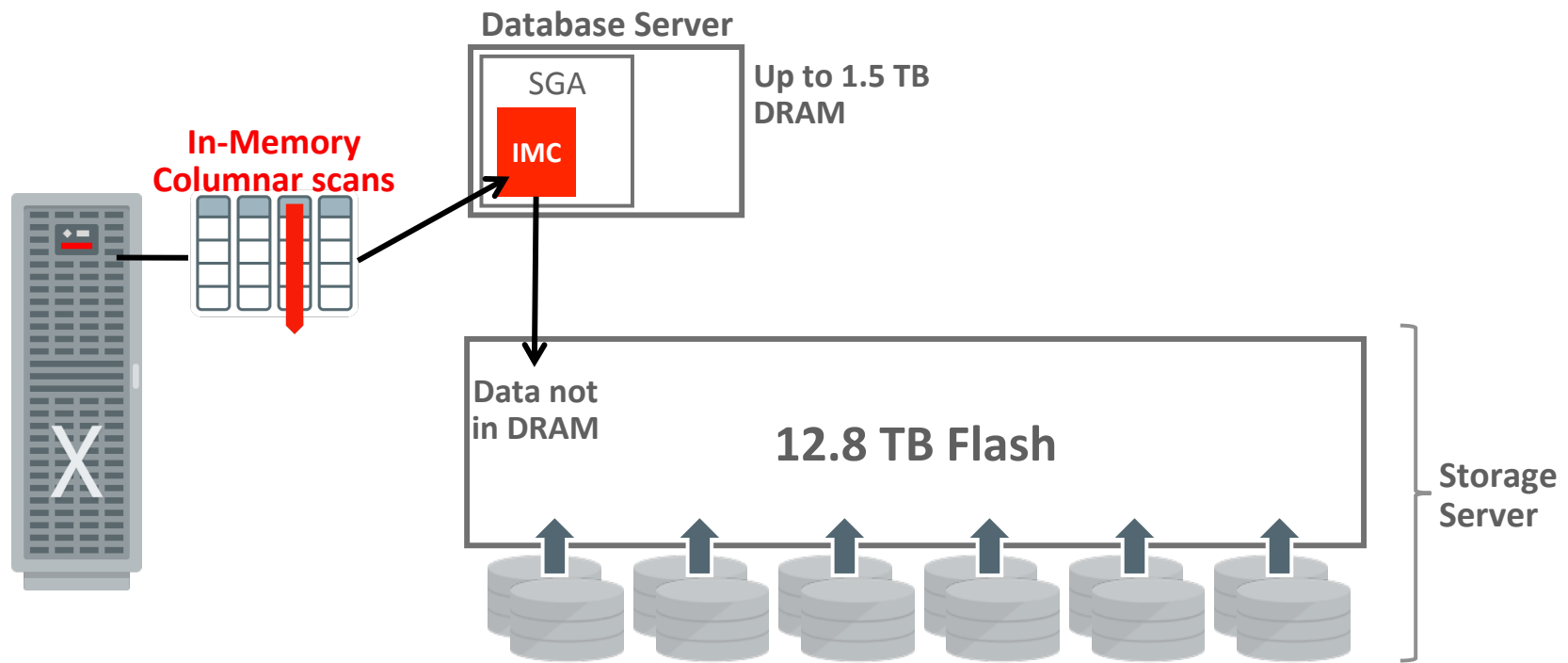**All-Flash Storage Array IO Path: many steps, each adds latency and creates bottlenecks**

Flash Chips → SSD Ctrl → **Array Heads** → Switches → Host HBA → CPU

Flash Chips     SAS/SATA     SAN/LAN     SAN/LAN     PCIe

# Redesigning Scan Offload for Memory Throughput

**NEW IN 12.2**

- With Exadata Flash throughput approaching memory throughput, SQL bottleneck moves from I/O to CPU

- Exadata will automatically transform table data into In-Memory DB columnar formats in Exadata flash cache
  - Dual format architecture extended from DRAM to flash

- Enables fast vector processing for storage server queries
  - Smart Scan results sent to DB using In-Memory Columnar format to reduce DB CPU usage

- Uniquely optimizes next generation flash as memory

**In-Memory Columnar scans**

**In-Flash Columnar scans**

ORACLE®

# In-Memory Columnar Formats in DRAM (pre 12.2.1.1.0)
## Super-Fast Scans from Memory, but All Queries Complete

**Database Server**

SGA

**IMC**

**Up to 1.5 TB DRAM**

**In-Memory Columnar scans**

**Data not in DRAM**

**12.8 TB Flash**

**Storage Server**

# In-Memory Columnar Formats in Flash Cache (12.2.1.1.0)

**3 - 4x Overall Analytics Performance Improvement**

**Database Server**

SGA

**IMC**

**Up to 1.5 TB DRAM**

**In-Memory Columnar scans**

**In-Flash Columnar**

*Extends In-Memory Column Store into Flash*

12.8 TB Flash x 3 **= 38.4 TB** (or more)

**IMC (In-Memory Columnar) data**

**Storage Server**

**Hybrid Columnar Compressed Data**

ORACLE

# Smart Analytics: Join and Aggregation Smart Scan

**NEW** DB In 12.2

- Extend In-Memory Aggregation technique into storage (vector joins and vector aggregation)

- Find Sales per country

```
SELECT /*+ VECTOR_TRANSFORM  */ country_id,
sum(amount_sold) amount_sold
FROM customers, sales
WHERE customers.cust_id = sales.cust_id
GROUP BY customers.country_id
ORDER BY customers.country_id;
```

- Storage cells scanning sales fact table return tuples

  `{country_id, sum_amount_sold }`

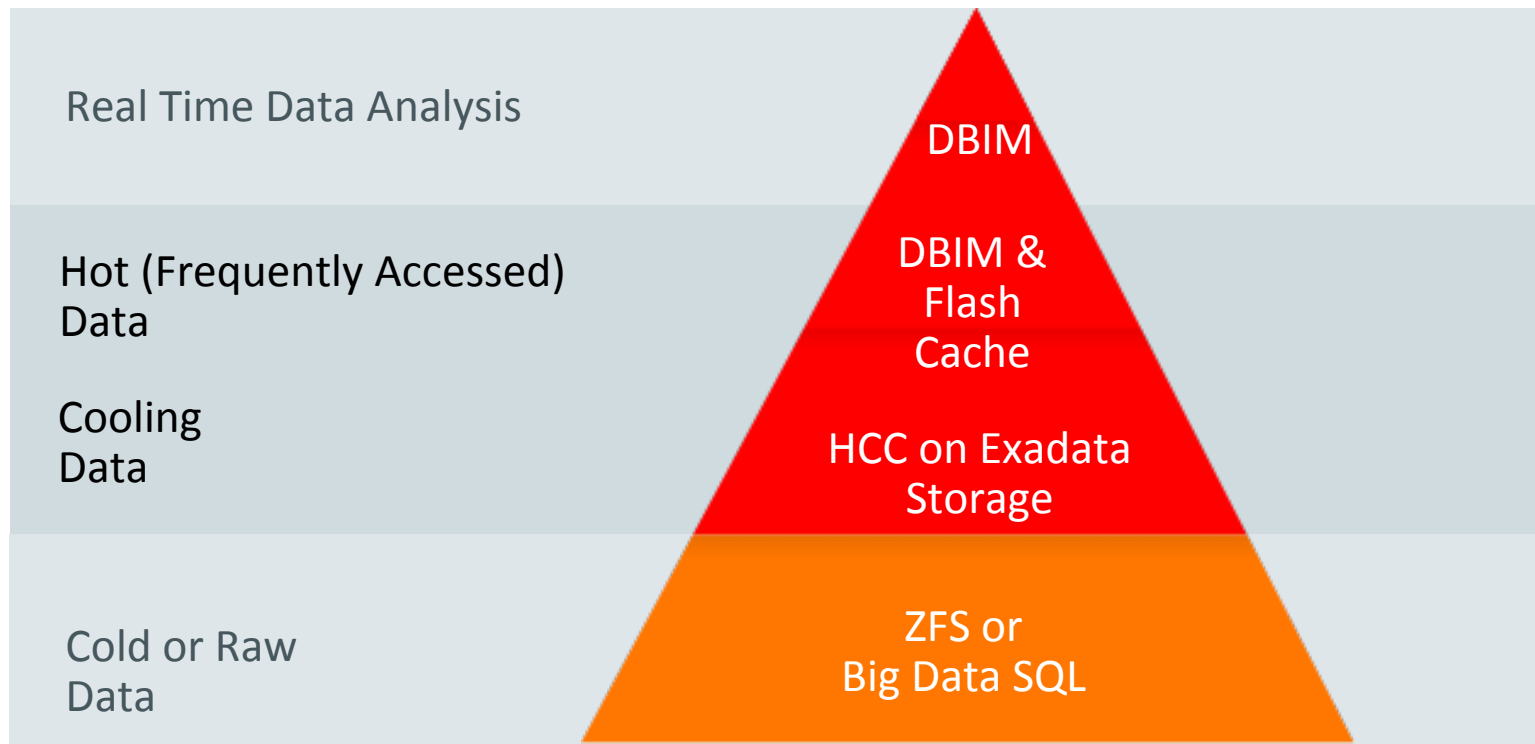- Join and Aggregation offloaded to the storage server

# **Smart Analytics:** More Smart Scan Enhancements

- Smart Scan enhancements for XML and JSON
  - JSON_EXISTS, JSON_VALUE, JSON_QUERY,
  "IS JSON" and "IS NOT JSON"
  - XML: XMLExists, XMLCast(XMLQuery())

- Significant speedup in JSON analytic workloads

```
select count(*)
  from pictures
  where json_value(photo, '$.tag')
  like '%spain%';
```
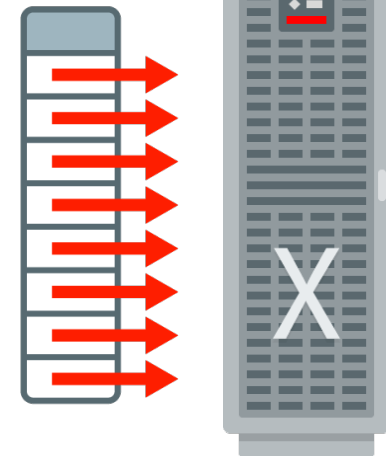


ORACLE®

# Data Tiering

| Real Time Data Analysis | DBIM |
|---|---|
| Hot (Frequently Accessed) Data | DBIM & Flash Cache |
| Cooling Data | HCC on Exadata Storage |
| Cold or Raw Data | ZFS or Big Data SQL |

ORACLE®

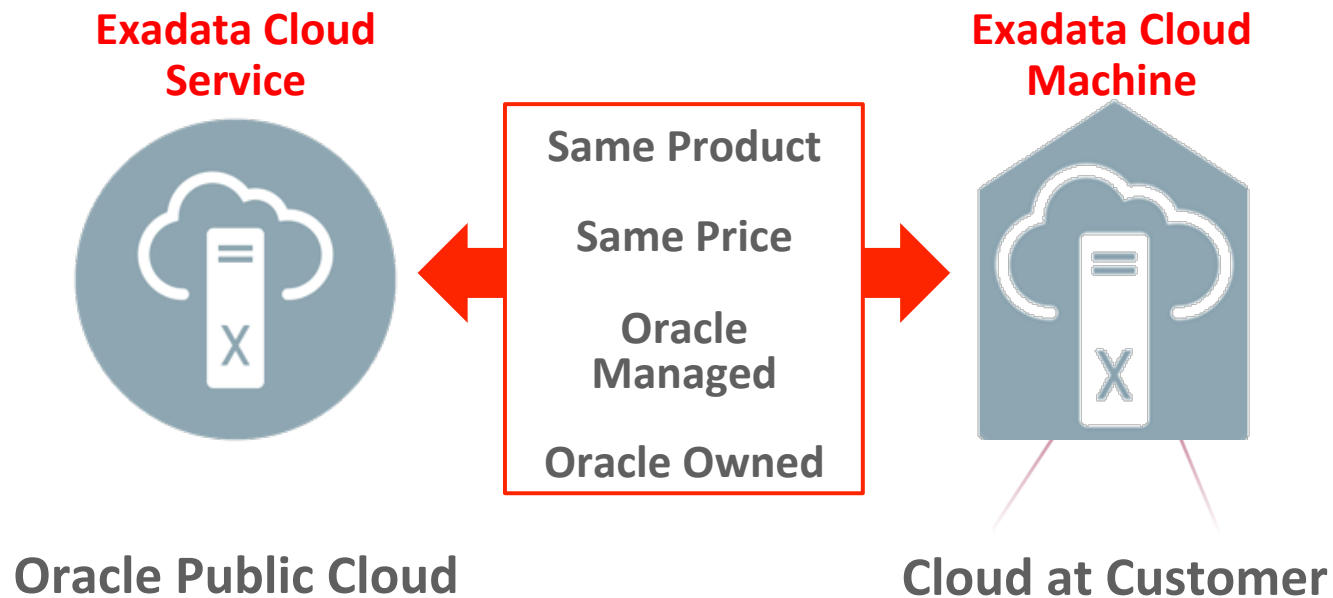# Smart Analytics: Smart Write Bursts and Temp IO in Flash Cache

- Write throughput of four flash cards has become greater than the write throughput of 12-disks

- When database write throughput exceeds throughput of disks, Smart Flash Cache intelligently caches writes

- When queries write a lot of temp IO, Smart Flash Cache intelligently caches temp IO
  - Writes to flash for temp spill reduces elapsed time
  - Reads from flash for temp reduces elapsed time further

- Smart Flash Cache prioritizes OLTP data and does not remove hot OLTP lines from the cache

- Smart flash wear management for large writes

- Supports Database 11.2.0.4, 12.1.0.2 and 12.2.0.1

**Write Bursts and Temp IO in Flash Cache**

**Accelerates Large Joins and Sorts and Large Data Loads**

ORACLE®

# Exadata Cloud – Your Way

**Exadata Cloud Service**

**Exadata Cloud Machine**

Same Product

Same Price

Oracle Managed

Oracle Owned

**Oracle Public Cloud**

**Cloud at Customer**

ORACLE®

# Exadata Customer Case Studies

# NTT docomo : MoBills（Mobile Billing System）

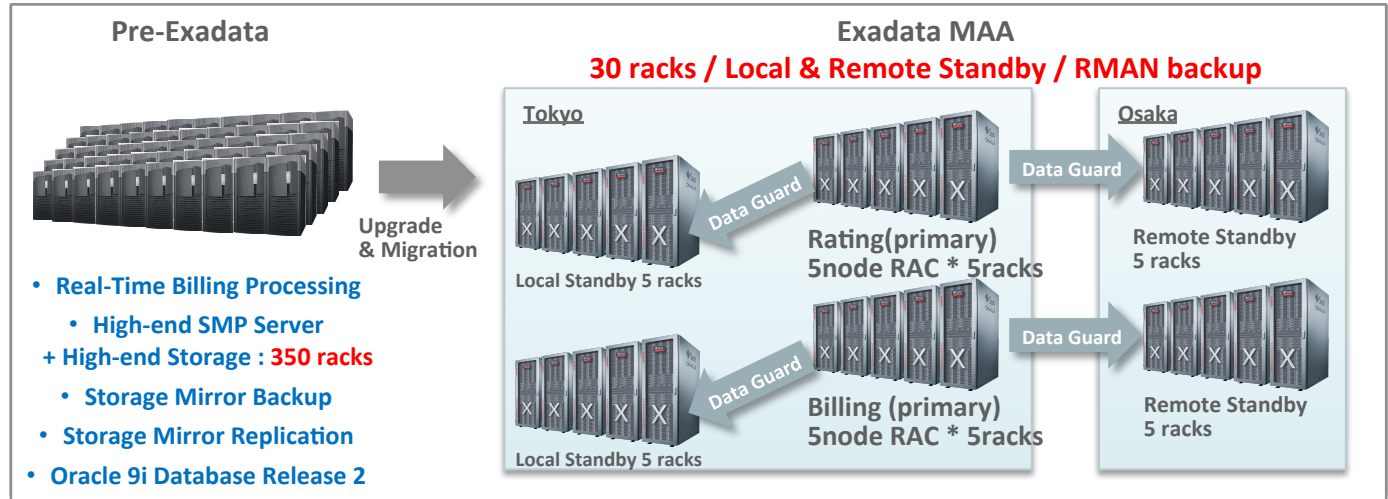| Benefits | Faster Billing Processing | Maximum Availability | Reduced Operational Cost | Reduced Introduction Cost | Data Center Cost Savings |
|---|---|---|---|---|---|
| *"MoBills is a very important position as a mission-critical system to promote efforts toward the realization of "+d". Oracle Exadata is running very stable as a expected performance. We will continue to use the "Oracle Exadata" and we would like to establish a further advantage for our business."* - Shimamura, Manager, Information System Department, NTT docomo | **10X speedup** 01010 010101 0110 **3 million SQL /sec** | **Local & Remote Standby** 24 7 | **50%** | **25%** | **90% Space Reduction** |

## Business Objectives

- Real-Time Billing Platform for 66 million customer
- Dramatically improve performance and availability
- Reduce cost and complexity

## Solution

- Oracle Exadata : 30 racks
- Oracle MAA (RAC / Active Data Guard - Local & Remote Standby database)

### Pre-Exadata

Upgrade & Migration

- **Real-Time Billing Processing**
  - **High-end SMP Server**
  - **+ High-end Storage : 350 racks**
  - **Storage Mirror Backup**
- **Storage Mirror Replication**
- **Oracle 9i Database Release 2**

### Exadata MAA
### 30 racks / Local & Remote Standby / RMAN backup

**Tokyo**

Local Standby 5 racks

Data Guard

**Rating(primary) 5node RAC * 5racks**

Data Guard

Local Standby 5 racks

Data Guard

**Billing (primary) 5node RAC * 5racks**

Data Guard

**Osaka**

Remote Standby 5 racks

Remote Standby 5 racks

ORACLE®

# DCM Holdings : System Consolidation of 3 companies

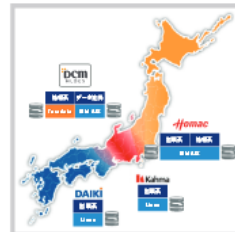| Benefits | Faster Batch Processing Reduced Introduction Cost | High Consolidation Ratio Improved Manageability | Simplified Support |
|---|---|---|---|
| *Realized the Database consolidation and integration due to the high performance provided by Oracle Exadata. And, Oracle Database 12c Multitenant Architecture also achieved high consolidation ratio while maintaining the independence of each group companies. Platinum Service could improve the service level, Oracle Full-stack products could provide One-Stop Support.* | **Standardization**<br><br>2X speedup     40% Off | **Multitenant Architecture**<br><br>6DBs Consolidation | **Oracle Full Stack**<br><br>Non Stop Support |

## Business Objectives

- $10 billion Sales, Faster M&A
- High Consolidation ratio and improve service level
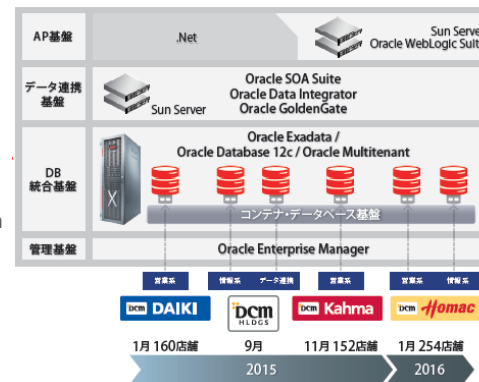- Reduce operational cost

## Solution

- Oracle Database 12c Multitenant on Exadata
- Oracle Full Stack (Middleware and Server products)

**Pre-Exadata**

Integration / Consolidation

**Oracle Multitenant on Exadata**

| AP基盤 | .Net | Sun Server Oracle WebLogic Suite |
| データ連携基盤 | Sun Server | Oracle SOA Suite Oracle Data Integrator Oracle GoldenGate |
| DB統合基盤 | X | Oracle Exadata / Oracle Database 12c / Oracle Multitenant<br>コンテナ・データベース基盤 |
| 管理基盤 | | Oracle Enterprise Manager |

DAIKI    DCM HLDGS    Kahma    Homac

1月160店舗    9月    11月152店舗    1月254店舗

2015    2016

- **Consolidation and Integration 3 group companies (Homac, Kahma and DAIKI) of system infrastructure**
- **Replaced from IBM p Servers**
- **Teradata Migration to Exadata**
- **Oracle Database 12c Multitenant**
- **Platinum Service**
- **Zero down time System Migration by using GoldenGate**

# Sprint: Call Data Record - Data Warehouse

| Benefits | Faster Queries | Faster Reports | Storage Savings | Maximum Availability | Data Center Cost Savings |
|---|---|---|---|---|---|
| | > 10x | 24 X | 6 x | No unplanned downtime | 3:1 Consolidation |

*"We reduced the queries from 30 seconds down to sub-second response time. Quick information, quick queries give Customer Care the ability to do their job better and meet the customer's needs."*
**- Richard Ewald, Senior Technical Architect, Data Warehousing**

15 billion transactions/day | 7 days to 7 hours | + removed 150 TB Indexes

## Business Objectives

- Improve performance
- Improve sustainability
- Improve availability and maintainability

## Solution

- Full Rack (Prod), Half Rack (Dev/Test); ZFS
- Storage Expansion
- Half Rack (Prod)

### Pre-Exadata

Sun Fire E6900

EMC / IBM / NetApp Storage

Sun M9000

- **4 x Sun Fire E6900, 1 x M9000**
- **Mixed Storage**
- **Multiple backup systems**
- **90 Day CDR DW 1.15 PB**
- **Oracle DB 11gR2**

### Production

2012 X3-2

2014 X4-2 Storage Expansion

2015 X5-2

ZFS Storage

- **Exadata X3-2 Full Rack**
- **HCC: 950 TB to 150 TB**
- **ZFS Storage Appliance (Backup)**
- **Exadata Storage Expansion**
- **Exadata X5-2 Half Rack**

### Dev/Test

Auto Service Request

Oracle Platinum Services

- **Exadata X3-2 Half Rack**

# Pulte Group: Multitenant Consolidation

PulteGroup™

| Benefits | | Business Impact | Faster Applications | Lower Admin & Support Costs | Cost Savings |
|---|---|---|---|---|---|
| *"Exadata delivered tremendous improvements in productivity. Users no longer have to wait for data. Data sharing is now real time."* <br> - Brian Pawlik, IS Manager, Pulte Homes | | 40% Productivity ⬆ Monthly Close 2 Days Faster | 2x -15x Faster | 40% Reduction | 40% CapEx |

## Business Objectives

- Scalability
- Supportability
- Sustainability

## Solution

- quarter rack & eighth rack

**Pre-Exadata**

IBM P7    EMC Storage

- IBM P7
- EMC storage arrays

**Exadata Quarter Rack**
Production / Standby / Test Dev / UAT

**Exadata Eighth Rack**
Disaster Recovery

Active Data Guard
WAN @ 800 miles

- Infor Lawson S3 ERP; Rebate Tracking
- Consolidate 35 DBs: 4 CDBs, 35 PDBs
- Production, Local Standby and QA
- Primary databases: > 5 TB

ORACLE®

# Exadata Advantages Increase Every Year

**Dramatically Better Platform for <u>All Database Workloads</u>**

*Smart Software*

*Smart Hardware*

- Exadata Cloud Machine
- Exadata Cloud Service
- In-Memory Columnar in Flash
- Smart Fusion Block Transfer
- In-Memory Fault Tolerance
- Direct-to-wire Protocol
- JSON and XML offload
- Instant failure detection
- Network Resource Management
- Multitenant Aware Resource Mgmt
- Prioritized File Recovery
- 3D V-NAND Flash
- IO Priorities
- Data Mining Offload
- Offload Decrypt on Scans
- Software-in-Silicon
- Tiered Disk/ Flash
- PCIe NVMe Flash
- Database Aware Flash Cache
- Storage Indexes
- Columnar Compression
- Smart Scan
- InfiniBand Scale-Out
- Unified InfiniBand
- DB Processors in Storage
- Scale-Out Storage
- Scale-Out Servers

# Integrated Cloud
## Applications & Platform Services

ORACLE®