# How to Create Massively Scalable Database Applications

Doug Hood

@ScalableDBDoug

Consulting Member of Technical Staff

Product Manager TimesTen In-Memory Database
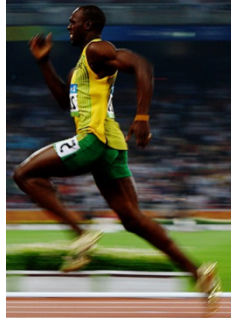
May 16, 2019

# Agenda

**1** Latency, Throughput and Scalability

**2** Scale-up vs Scale-out

**3** Scale-out Architectures

**4** Trivial Scalability Benchmarks

**5** Scaling a Customer Workload

**6** Summary and Q & A

**ORACLE**®

# Latency, Throughput and Scalability

| | | |
|---|---|---|
| **Latency** |  | **How quickly can one operation complete**<br><br>**One sprinter in 9.58 seconds**<br>**~ 40 km/h for 100M [2009]** |
| **Throughput** |  | **How quickly can many operations complete**<br><br>**Ten sprinters in under 11 seconds**<br>**~ 40 km/h for 100m [2009]** |
| **Scalability** |  | **By adding more resources can throughput keep increasing**<br><br>**33 cars on 2.5 mile oval track**<br>**~250 km/h for 804 km [Indy 500, 2017]** |

# You can only go so big



IBM z14
z14 Microprocessor
32 CPUs @ 5.2 GHz
- 10 cores, 20 threads
**8 TB DRAM**



Oracle Super Cluster
SPARC M8
16 CPUs @ 5.1 GHz
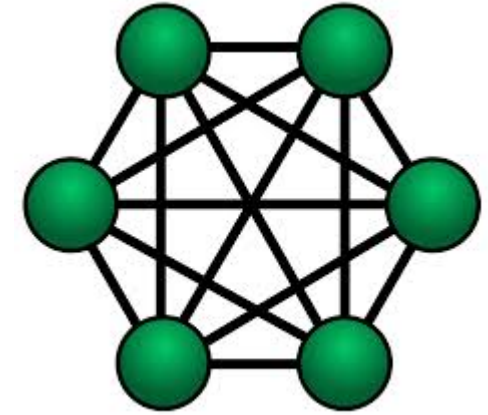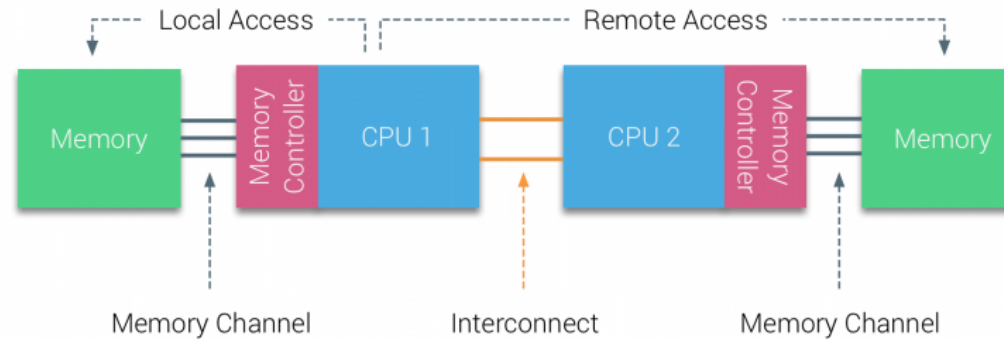- 32 cores, 256 threads
**16 TB DRAM**



HPE Superdome Flex
Intel Xeon
32 CPUs @ 3.6 GHz
- 28 cores, 56 threads
**48 TB DRAM**



SGI Altix 4700
Intel Itanium 2
2048 CPUs @ 900 MHz
- 2 cores, 4 threads
**128 TB DRAM**

# Vertical Scaling Limits

- Only so many CPUs interconnected
- NUMA limits
- Complexity & Cost
- Niche Market

Local Access — Remote Access

Memory | Memory Controller | CPU 1 | CPU 2 | Memory Controller | Memory

Memory Channel     Interconnect     Memory Channel

8+ Sockets

4-8 Sockets

1-2 Sockets

# Horizontal Scaling hardware

- Use **cheap/fast Linux x8664 servers**, eg Oracle Sun X7-2

- NUMA affects are minimal

- Commodity servers keep getting **faster**, **cheaper** and **more powerful**

- 1.5 TB DRAM [Persistent Ram coming, **Intel/Oracle PMem demo**]

- Two Intel Xeon 8164 @ 2.2 GHz, 26 cores

- Up to eight NVMe SSDs

- **42 1U servers per Rack:**
  - 2*42 = **84 CPUs**
  - 1.5 * 42 = **63 TB RAM**

# Lower Latency with TimesTen Cache

| Query | Oracle | Cache |
|-------|--------|-------|
| Q1    | 43     | 3     |
| Q2    | 69     | 6     |
| Q3    | 105    | 8     |
| Q4    | 121    | 20    |
| Q5    | 140    | 18    |
| Q6    | 163    | 19    |
| Q7    | 231    | 18    |

**Oracle 11.2.0.4 RAC**
**RAC nodes were Oracle Sun X7-2L**
**NVMe Storage**
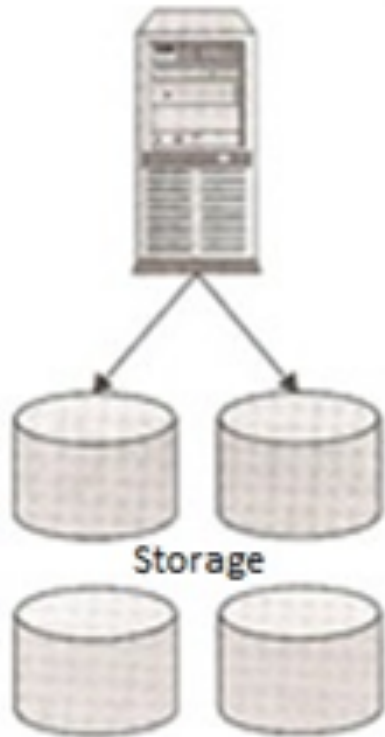**Over 50 Million Users**

**Application Tier Database Cache (TimesTen)**
**Ran on the same nodes as the production RAC**
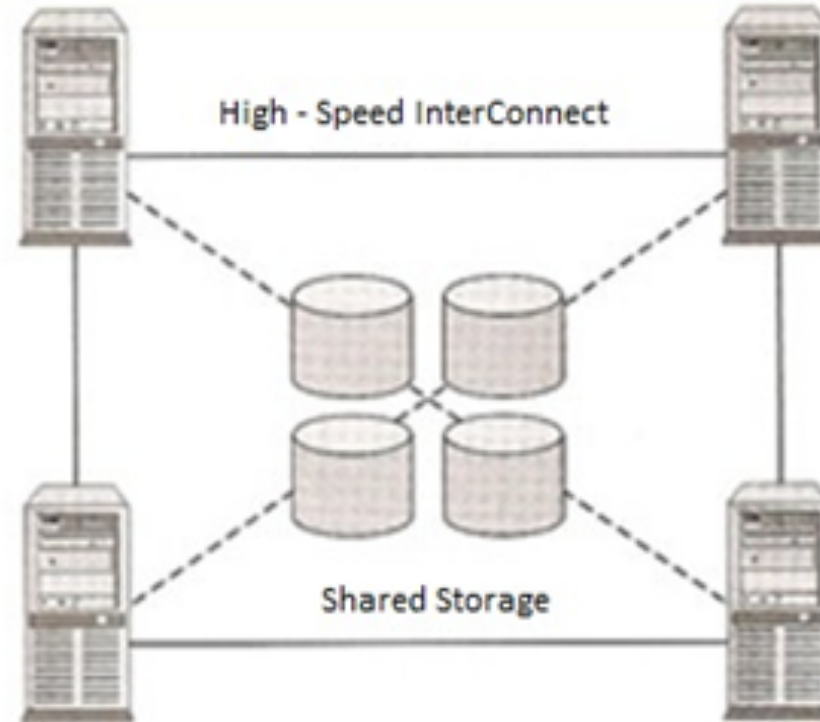**5 table joins for 100s of millions of rows of data**

# Latency is in Micro Seconds ...

# Oracle Database & Real Application Clusters Architecture


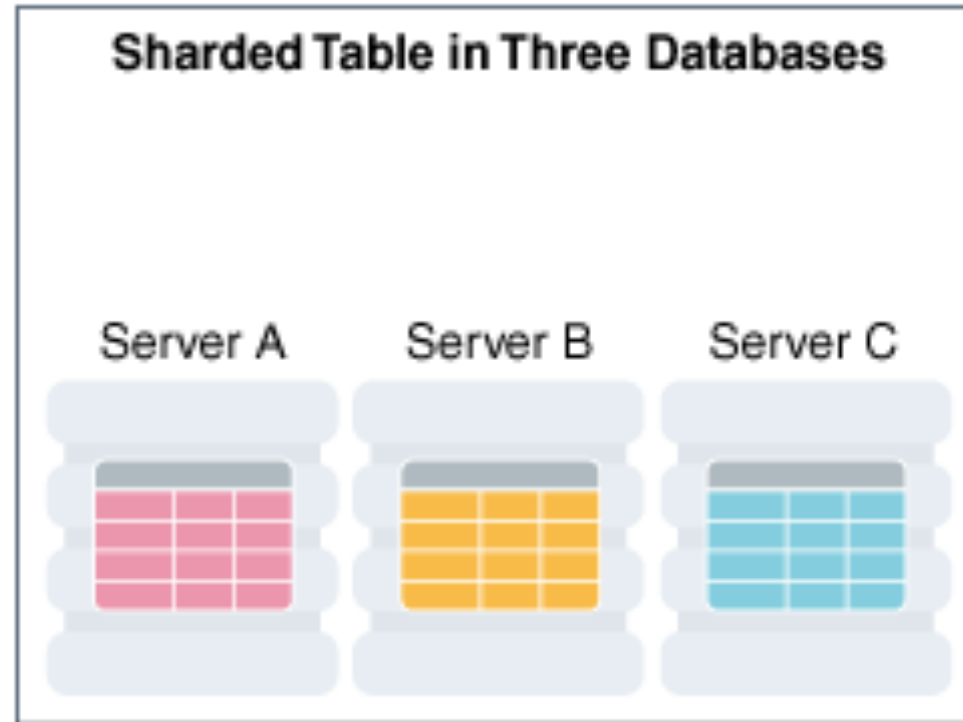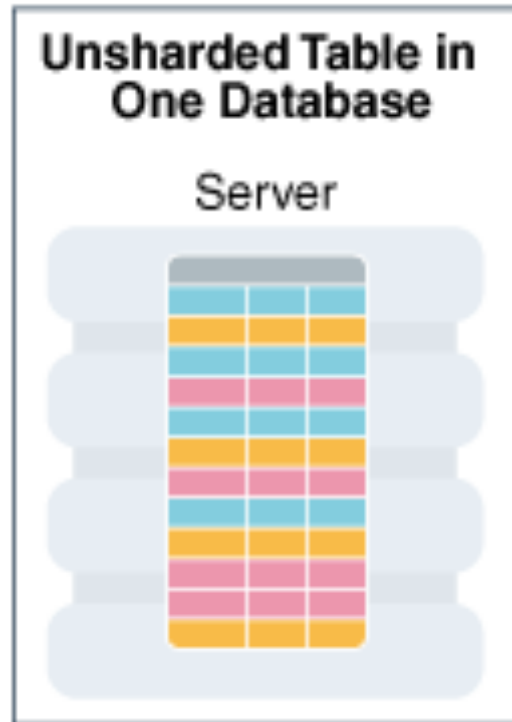
**Oracle Database**
- Single Instance
- Single DB image

**Oracle Real Application Clusters**
- Multiple Database Instances
- Single DB image
- Shared Storage

**Oracle Exadata**
- Multiple Database Instances
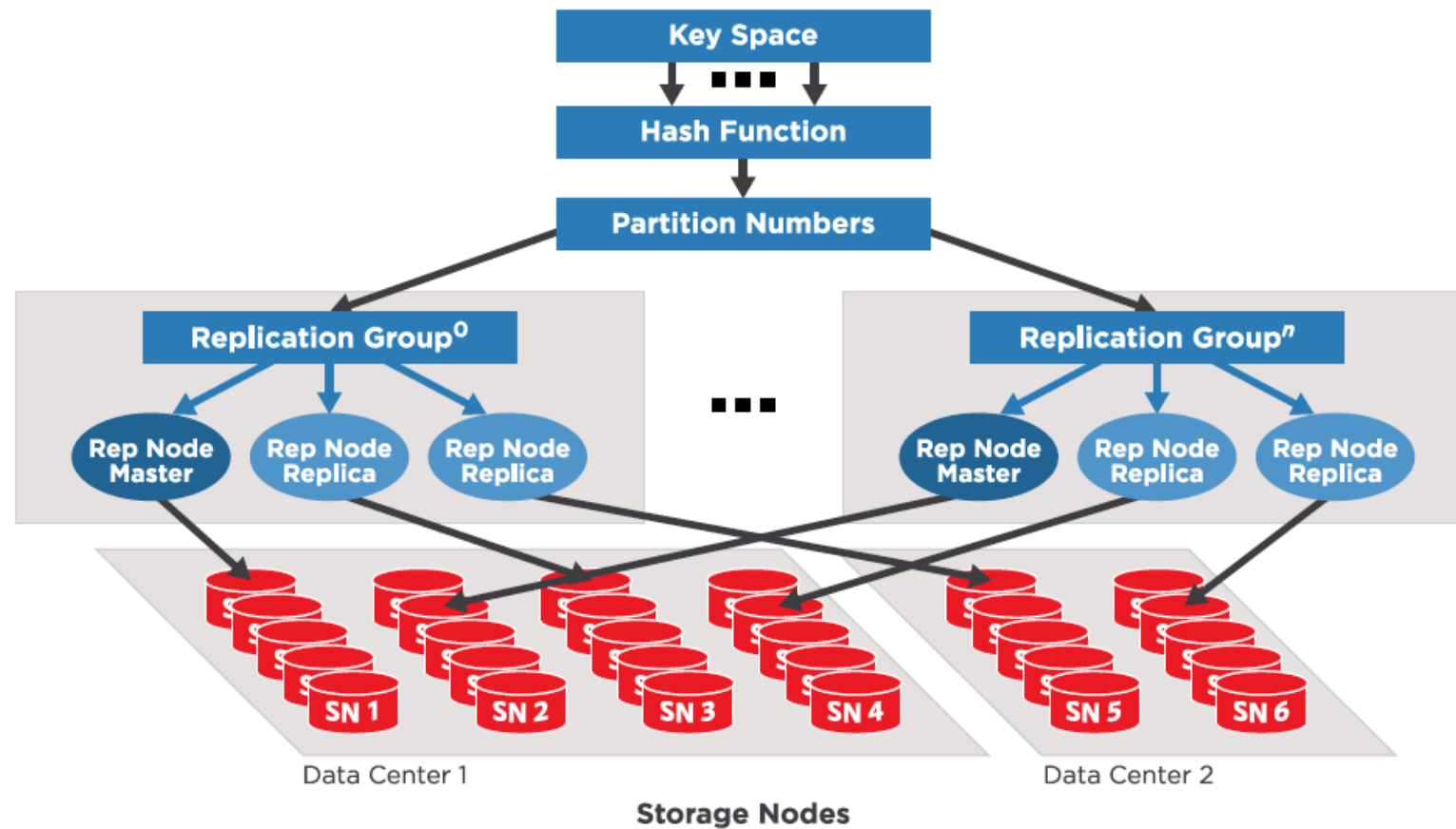- Single DB image
- Shared Storage

# Oracle Sharding Architecture

**Oracle Sharding**
- Multiple Database Instances
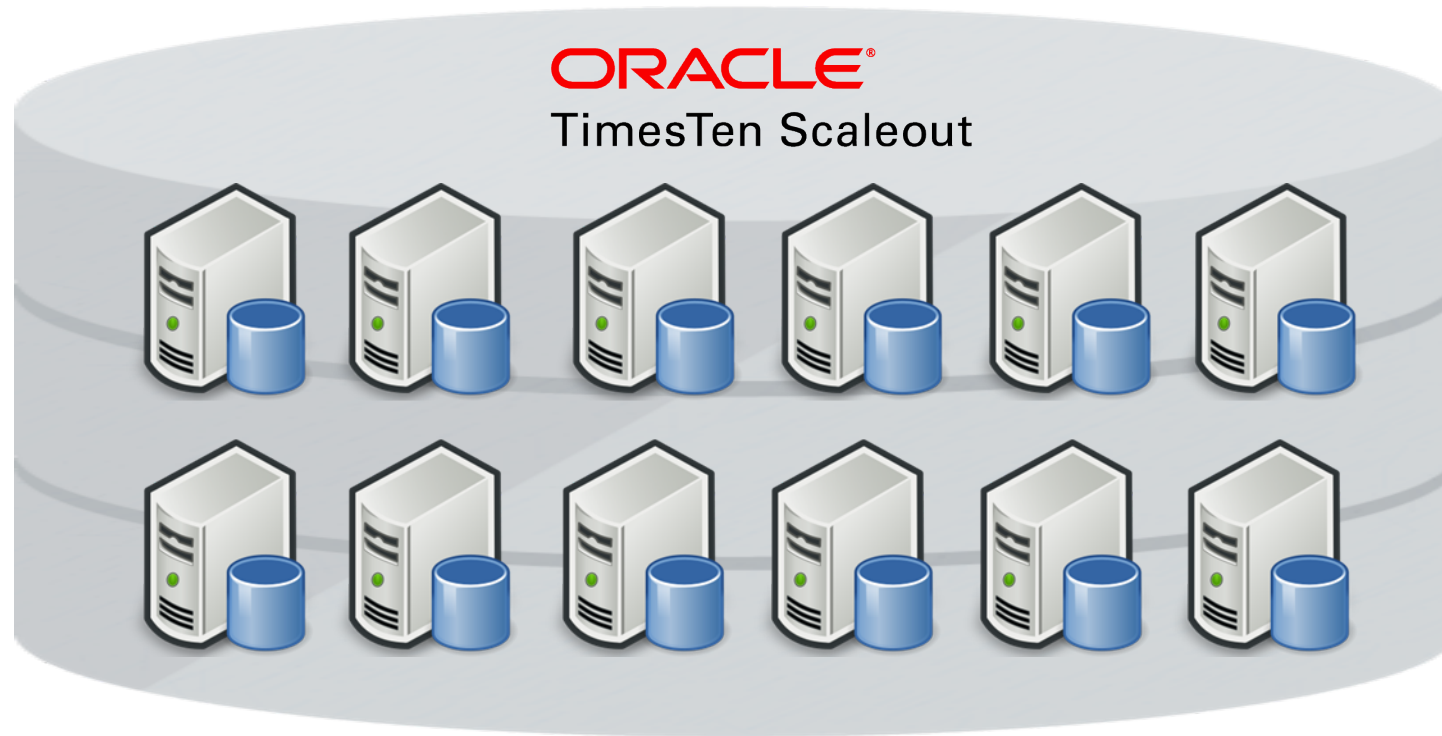- Multiple DB images
- Independent Storage

# Oracle NoSQL Architecture



**Oracle NoSQL**
- Multiple 'DB' Instances
- One DB image
- Independent Storage

# Oracle TimesTen Scaleout Architecture



**Oracle TimesTen Scaleout Architecture**
- Multiple Database Instances
- Single DB image
- Shared Nothing

# Summary of how to Scale Database Apps



12

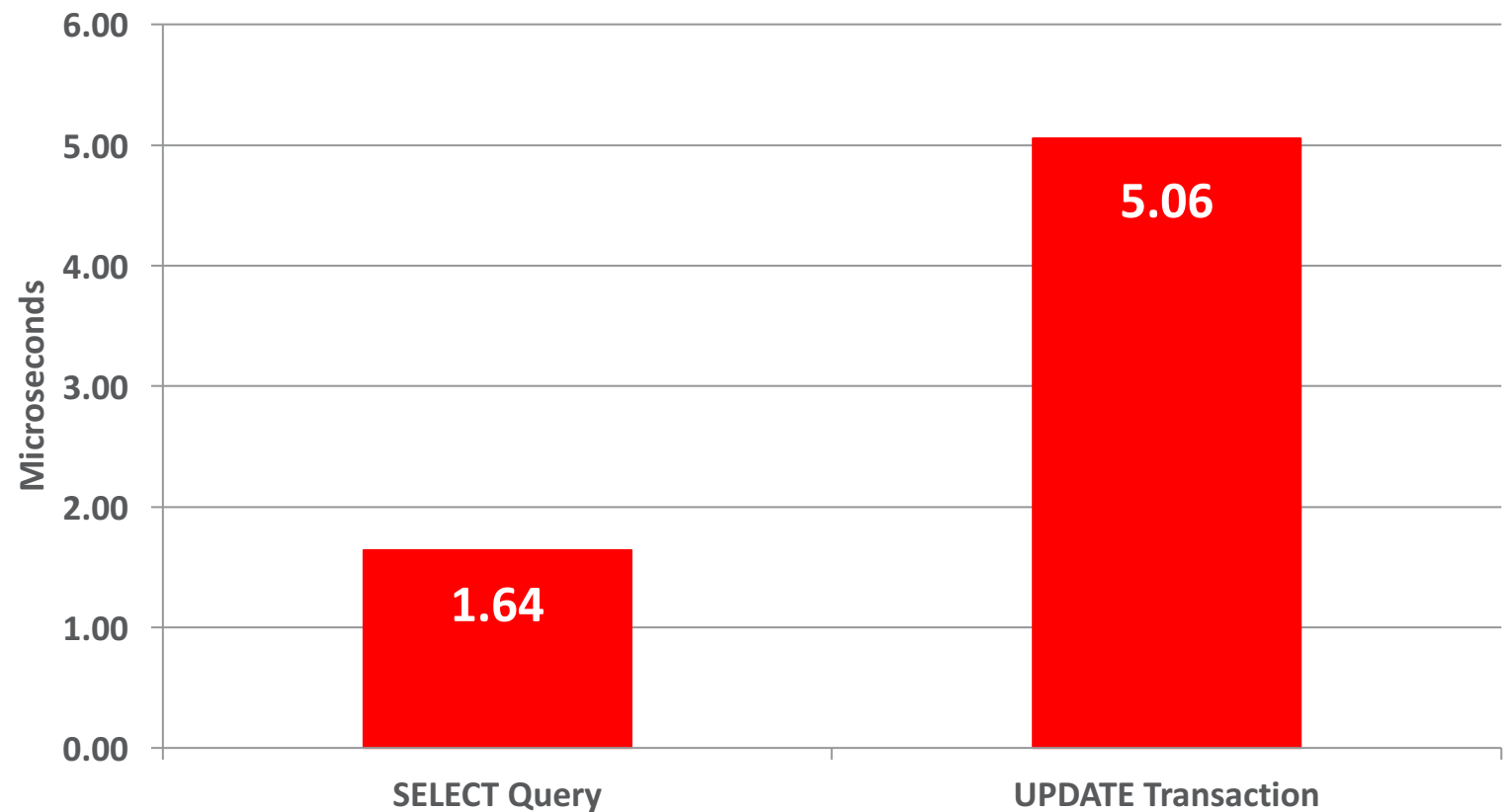# Summary of how to Scale Database Apps

- Do not do dumb things

- Tune your SQL

- Use PLSQL stored procedures intelligently

- Use good hardware

- Scale-up with Sun SuperCluster

- Scale-out with Exadata

- Scale-out with Application Tier Database Cache or TimesTen Scaleout

# Low Latency - **Microseconds** Response Time

*select directory_nb,*
    *last_calling_party,*
    *descr*
*from vpn_users*
*where vpn_id = :1*
*and vpn_nb= :2*

**TPTBM Read and Update
E5-2699 v4 @ 2.20GHz
2 socket, 22 cores/socket,
2 threads/core
TimesTen 11.2.2.8.0
(100M rows, 17GB data)**



Bar chart — Microseconds (y-axis 0.00 to 6.00):
- SELECT Query: 1.64
- UPDATE Transaction: 5.06

# Some Throughput & Scalability Benchmarks

- YCSB : **Y**ahoo **C**loud **S**erving **B**enchmark
  - Developed at Yahoo for Cloud Scale workloads
  - Widely used to compare scale-out databases, NoSQL databases, and (non-durable) in-memory data grids

- A series of workload types are defined:
  - Workload A: 50% reads, 50% Updates
  - Workload B: 95% reads, 5% Updates
  - Workload C: 100% reads

- The YCSB Client cannot be changed
  - DB Vendors implement the DB Client interface in Java
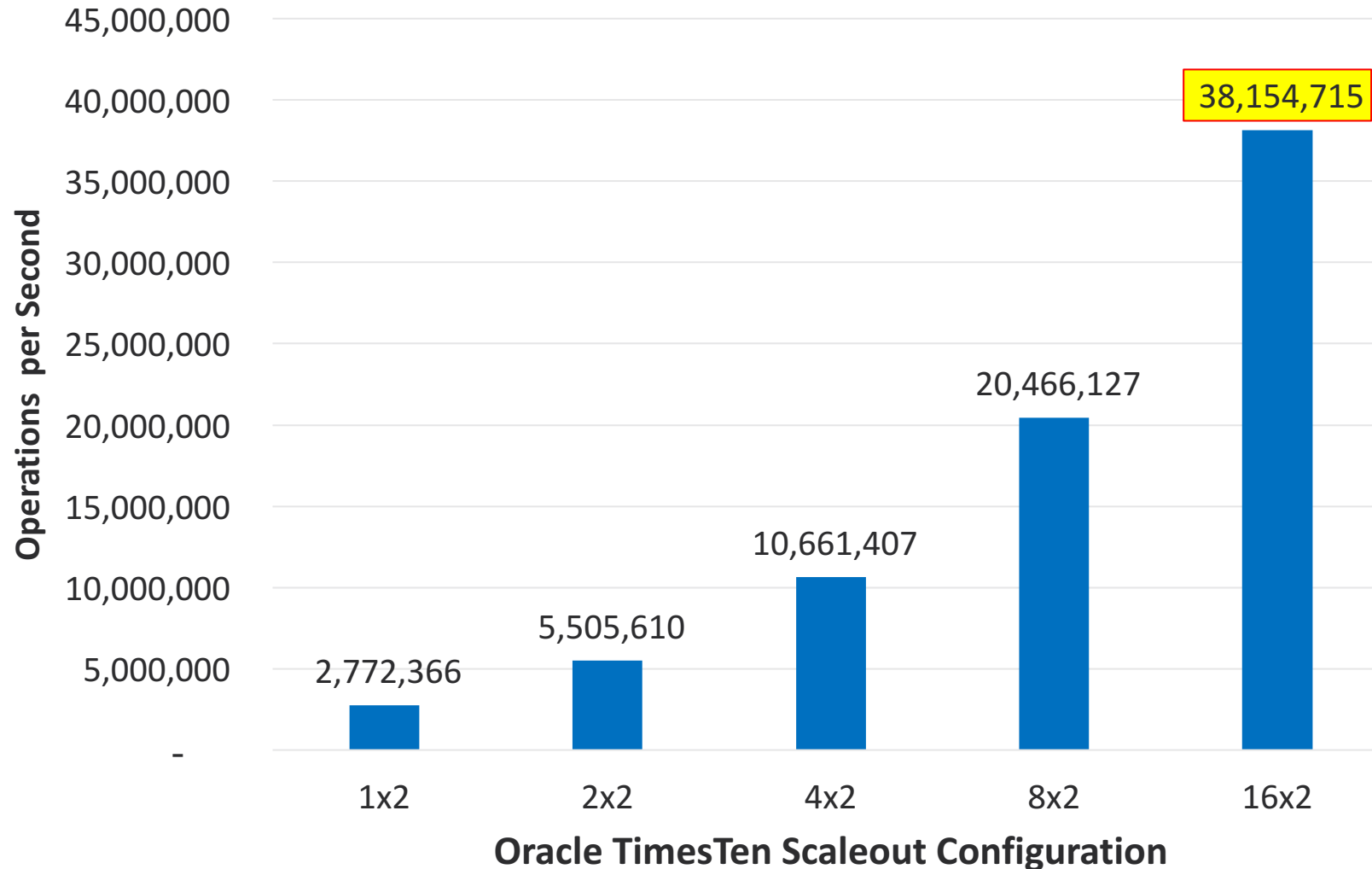  - The version and exact configuration matters

## Surveyed YCSB (Workload B) Results*

| Product | Type | Nodes | Ops/Sec |
|---|---|---|---|
| cassandra | NoSQL DB | 32 | **227 K** |
| mongoDB | NoSQL DB | 2 | **275 K** |
| SCYLLA | NoSQL DB | 3 | **715 K** |
| VOLTDB | Scale-Out RDBMS | 6 | **1.6 M** |
| AEROSPIKE | NoSQL DB | 8 | **1.6 M** |

*There is no official repository of YCSB results*
*These were the largest results we found online*

# TPTBM 80% Read 20% Update: **153 Million** Transactions/Sec



**TPTBM Configuration**

- 128-byte record
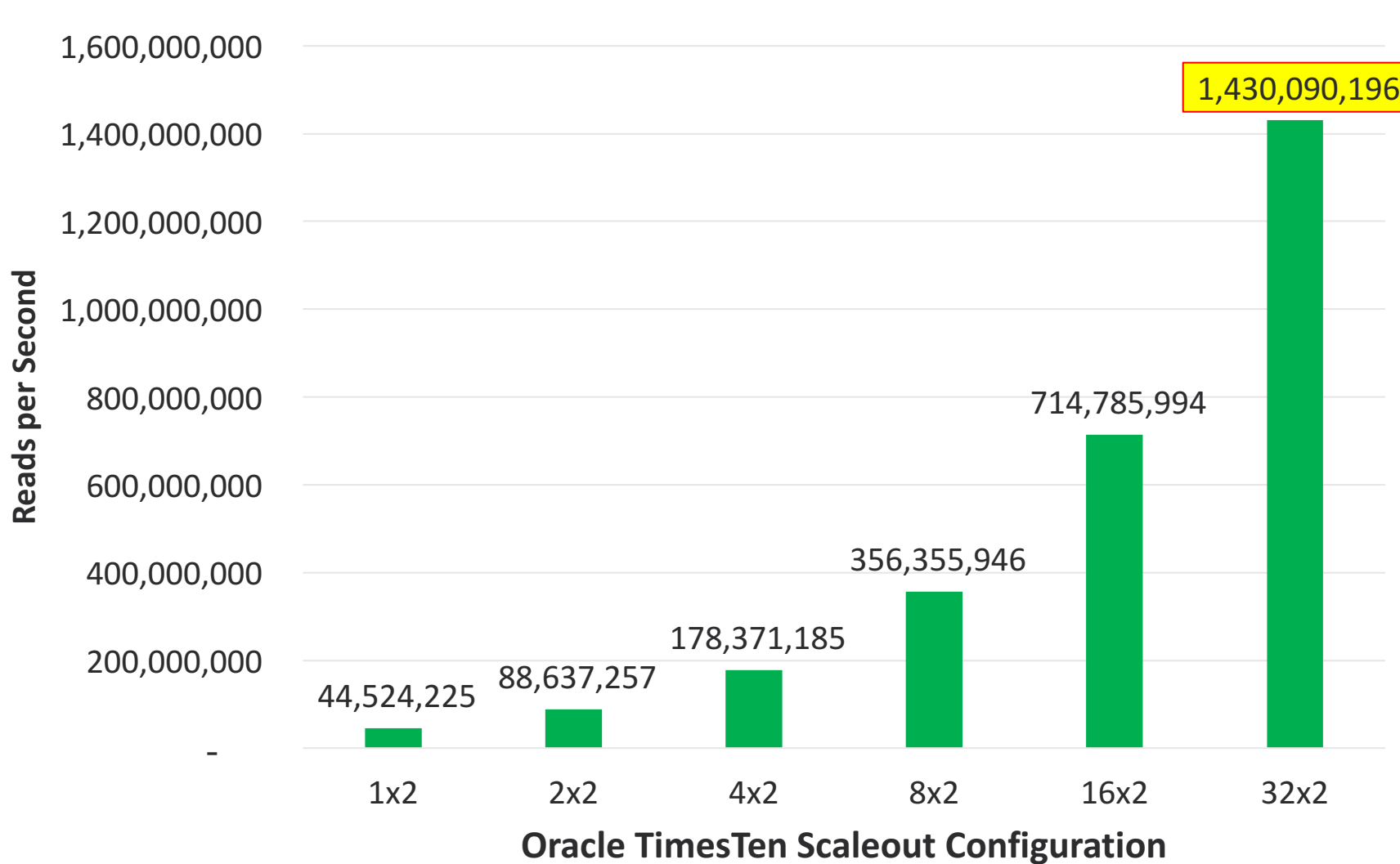- 100M records / Replica Set
- Uniform Distribution

**TimesTen Scaleout**

- 1 to 64 replica sets
- 1 replica per replica set

**Oracle Cloud Infrastructure**

- 32 * BM.DenseIO2.52
- Two TimesTen instances per compute node

Chart values by Oracle TimesTen Scaleout Configuration:
- 2x1: 5,695,071
- 4x1: 11,251,034
- 8x1: 22,611,122
- 16x1: 41,633,465
- 32x1: 81,746,166
- 64x1: 153,140,347

Y-axis: Transactions per Second
X-axis: Oracle TimesTen Scaleout Configuration

ORACLE®

# TPTBM 100% Read: **1.4 Billion Reads** Per Second!!



**Reads per Second** (y-axis)

| Configuration | Reads per Second |
|---|---|
| 1x2 | 44,524,225 |
| 2x2 | 88,637,257 |
| 4x2 | 178,371,185 |
| 8x2 | 356,355,946 |
| 16x2 | 714,785,994 |
| 32x2 | 1,430,090,196 |

**Oracle TimesTen Scaleout Configuration** (x-axis)

### TPTBM Configuration

- 128-byte record
- 100M records / Replica Set
- Uniform Distribution

### TimesTen Scaleout

- 1 to 32 replica sets
- 2 synchronous replicas per replica set

### Oracle Cloud Infrastructure

- 32 * BM.DenseIO2.52
- Two TimesTen instances per compute node

18

# What Hardware was Used?

**Oracle Sun X7-2**
- Two Intel Xeon 8164 @ 2.2 GHz, 26 cores
- 768 GB RAM
- Four NVMe SSDs
- Two 10G Ethernet
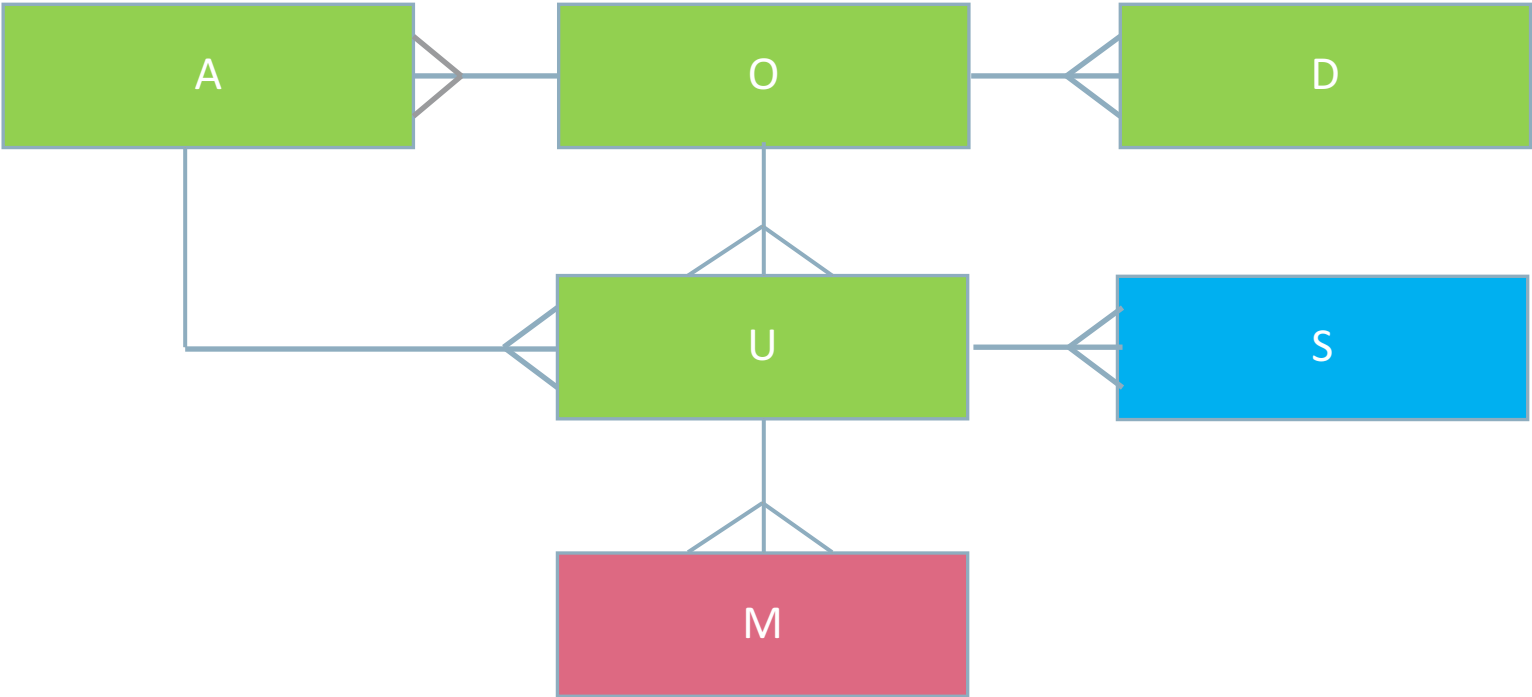
**Oracle Cloud Infrastructure**
- 32 * BM.DenseIO2.52

**ORACLE** CLOUD

# ORACLE®
# TimesTen Scaleout

## World's Fastest OLTP Database

# Subset of Customer's Data Model



**+ seven other tables for the 'write' workload**

# Critical Query

```
SELECT a.usr_id, …
FROM u, d, o, a
WHERE u.login_name = :loginName
AND u.dom_id = a.dom_id
AND u.usr_org_id = o.org_id
AND u.account_id = a.acct_id (+)
AND u.status <> :x;
```

```
SELECT mn_usr_id
FROM m
WHERE  mn_usr_id = uid
AND status = :y;
```

```
SELECT s.attr_name
FROM s
WHERE  s.entity_id = muid
AND (s.context = :p or b.context = :q)
AND (s.spid = :m or
       s.spid = :n or
       s.sid = :o)
ORDER BY b.attr_name;
```

# Critical Update Transaction

```
select something
from R1
where col1  = :x
and col2  = :y;


update  R1
set something  = :s
where col1  = :x
and col2  = :y;
```

```
select something
from R2
where col1  = :x
and col2  = :y;


update  R2
set something  = :s
where col1  = :x
and col2  = :y;
```

```
select something
from R3
where col1  = :x
and col2  = :y;


update  R3
set something  = :s
where col1  = :x
and col2  = :y;
```

```
select something
from R7
where col1  = :x
and col2  = :y;


update  R7
set something  = :s
where col1  = :x
and col2  = :y;
```

# Scale Up or Scale Out?



Four 5.1 GHz SPARC CPUs
256 hardware threads per CPU socket
64 MB L3 Cache
16 TB RAM
8 NVMe SSD for DB Storage + 12 Disks
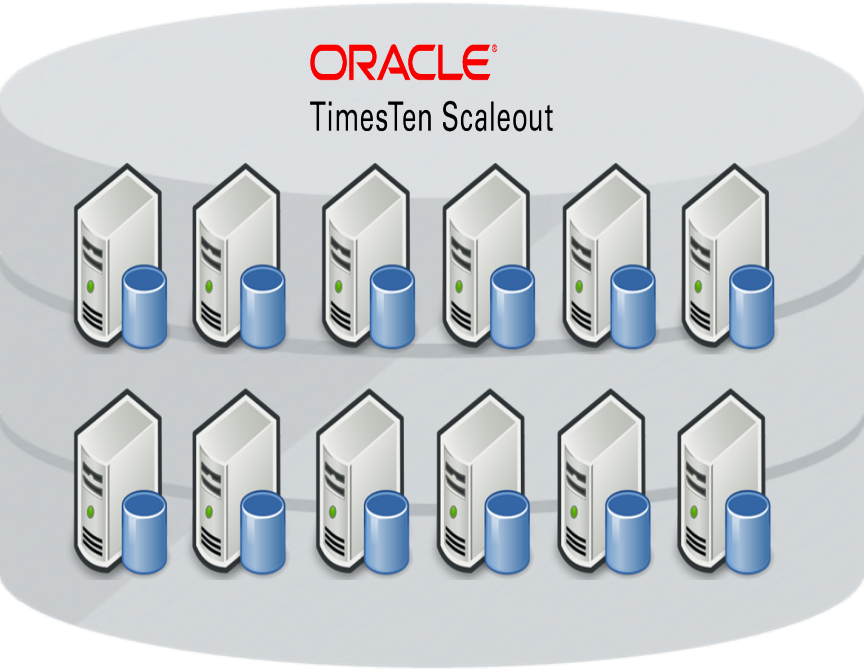40 G Infiniband
4 Quad 10G Ethernet

Oracle Database 11g



32 Core VMs
64 GB RAM
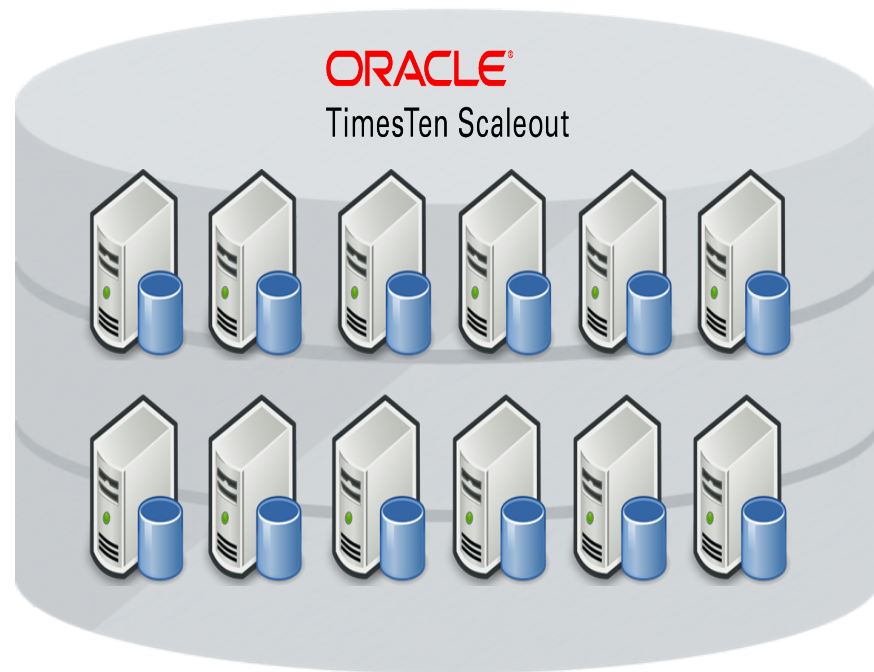Cinder Storage

# Best Case Architecture for customer workload



ORACLE®
TimesTen Scaleout

| Attribute | Value |
|---|---|
| Data Reads & Writes | 100% Local RAM |
| Storage Reads & Writes | 100% Local NVMe SSD |
| Storage Bottleneck | No |
| Fast CPU | Xeon |
| Number of CPU cores | 24 |
| Sufficient Memory | Yes. 320 GB |
| DB Tuned | Yes |
| App tuned | **No. Python without SQL prepares or binds** |

# Result: 11 Million Transactions / second

ORACLE®

# Worst Case Architecture for customer workload

TimesTen Scaleout

| Attribute | Value |
|---|---|
| Data Reads & Writes | 90% on a remote VM |
| Storage Reads & Writes | 100% remote [Cinder/Netapp] |
| Storage Bottleneck | Maybe. Network bound |
| Fast CPU | Xeon |
| Number of CPU cores | 32 |
| Sufficient Memory | No. Only 32 GB |
| DB Tuned | Yes |
| App tuned | Yes. ODBC with SQL prepares and binds |

## Result: 304K Transactions / second

# Some Results



**240K TPS**
60/40 Workload
IO Bound
ACID 1PC

**4 Socket SMP**



**< 168K TPS**
60/40 Workload
Network Bound
Eventual Cons

**Negative Scaling**



**168K TPS**
60/40 Workload
Network Bound
Eventual Cons

**37 Node Cluster**


TimesTen Scaleout

**304K TPS**
60/40 Workload
Network Bound
ACID 2PC

**10 Node Cluster**

# How Many Client Server SQL Network Round Trips ?

1. Select * from table where PK = :value;
2. Select * from table where PK between 10 and 20;
3. Update table set column = :X where PK = :value;
4. Update table set column = :X where PK between 1000 and 2000;
5. Select * from a, b, c, d where {non Cartesian Product}

A. One
B. Two
C. Three
D. Lots
E. It Depends

*How many server side network messages When tables are hash distributed?*

**ORACLE**®

# Data Distribution Methods
## Distribute Table Data by Hash, Reference or Duplicate

- Distribute by **Hash**
  - Primary key or user-specified columns
  - Consistent hash algorithm
  - Examples: Customers, Subscribers, Accounts

- Distribute by **Reference**
  - Co-locate related data to optimize joins
  - Based on FK relationship
  - Supports multi-level hierarchy

- Distribute by **Duplicate**
  - Identical copies on all elements
  - Useful for reference tables
  - Read and join optimization

Distribute by Hash → Customer

Distribute by Reference → Order

Distribute by Duplicate → Products

| Element 1 | | Element 2 | | Element 3 | | Element 4 | |
|---|---|---|---|---|---|---|---|
| 0 | David | 1 | Bill | 2 | Olaf | 3 | Chi |
| 4 | Igor | 5 | Sam | 6 | Henri | 7 | Simon |
| 8 | Tim | 9 | Charles | 10 | Jie | 11 | Chris |

| Element 1 | | | Element 2 | | | Element 3 | | | Element 4 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 16/6/15 | 2 | 5 | 16/2/22 | 5 | 6 | 16/5/10 | 3 | 3 | 16/3/1 |
| 6 | 8 | 16/3/22 | | | | | | | 4 | 11 | 16/2/5 |

| Element 1 | | Element 2 | | Element 3 | | Element 4 | |
|---|---|---|---|---|---|---|---|
| phone | 100 | phone | 100 | phone | 100 | phone | 100 |
| tablet | 200 | tablet | 200 | tablet | 200 | tablet | 200 |
| watch | 300 | watch | 300 | watch | 300 | watch | 300 |

# Scalability Challenges

- Four table joins with hash distribution for 'read workload' with (+)

- Seven queries + seven updates for 'write workload'

- Client Server round trips

- Not enough RAM [64 GB] per VM

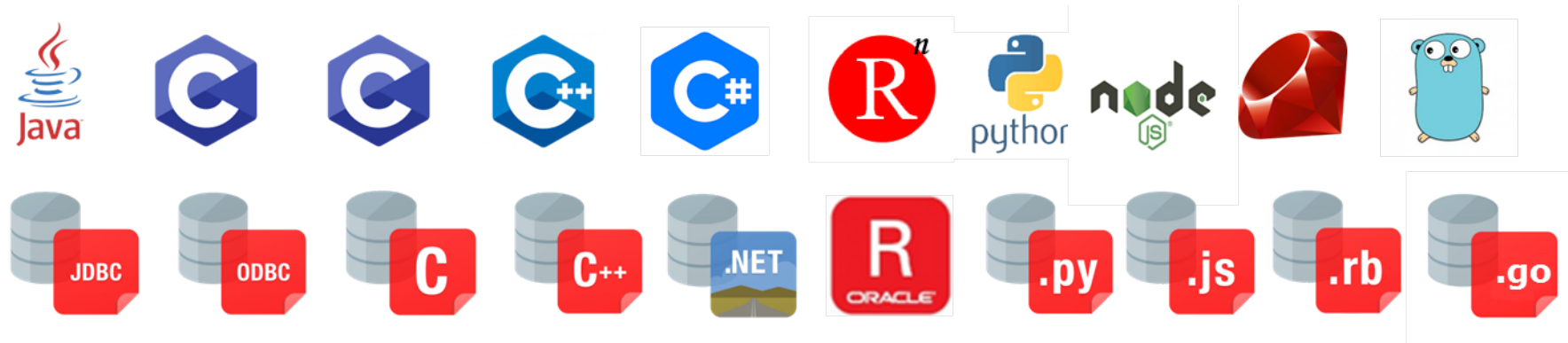- KVM + OpenStack Neutron networking overhead

# Techniques which helped Scalability

- Determine the best distribution clauses
  - The Distribution Advisor eliminates the guess work

- Determine the best indexes
  - The Index Advisor eliminates the guess work

- Prepare and Bind the SQL statements

- Check the explain plans

- Use Stored Procedures for the 'read' and 'write' transactions
  - Execute many statement in a single network round trip. Procedural logic + commit/rollback

- Use the Routing API
  - Determine where the data is to avoid network hops

- Use more DB nodes
  - The VNIC became network bound [ksoftirq]
  - Use more modes to lessen the load per VNIC

**TODO**
- TCP tuning
- RDMA

# TimesTen Scaleout SQL APIs

| API | Comment |
| --- | --- |
| JDBC | The same (JDBC 4.3) |
| ODBC | The same (ODBC 3.5.2) |
| OCI | The same (OCI 11.2.0.4.+) |
| R-Oracle | The same (OCI 11.2.0.4.+) |
| ODP.Net | The same (OCI 11.2.0.4.+) |
| PL/SQL | The same (11.2.0.4.+) |
| Python | The same (cx_Oracle, ODPI-C) |
| Ruby | The same (Ruby-ODPI, ODPI-C) |
| GoLang | The same (go-goracle, ODPI-C) |

ORACLE®

# TimesTen in On Premises

- TimesTen Scaleout requires :
  - Linux x8664 (glibc 2.12+)
    - Oracle Linux / Red Hat / CentOS 6.4+, 7+
    - Ubuntu 14.04+
    - SuSE 12+
  - JDK 8+
  - TCP/IP or IPoIB
  - A file system [eg ext4, not ext2 or ext3]
  - Enough RAM for the DB

# TimesTen Scaleout on OCI, AWS, Azure, Google

# Centralized Installation and Management

- All TimesTen Scaleout management and admin operations are performed from a single host
  - Installing software
  - Patching software
  - Configuration
  - Database creation and management
  - Backup and restore
  - Monitoring
  - Collecting diagnostics

- Command line interface
- SQL Developer (GUI) interface

```
-- Database is in Oracle type mode
create table APPUSER.ACCOUNTS (
        ACCOUNT_ID          NUMBER(10) NOT NULL,
        PHONE               VARCHAR2(16 BYTE) INLINE NOT NULL,
        ACCOUNT_TYPE        CHAR(1 BYTE) NOT NULL,
        STATUS              NUMBER(2) NOT NULL,
        CURRENT_BALANCE     NUMBER(10,2) NOT NULL,
        PREV_BALANCE        NUMBER(10,2) NOT NULL,
        DATE_CREATED        DATE NOT NULL,
        CUST_ID             NUMBER(10) NOT NULL,
    primary key (ACCOUNT_ID),
    constraint FK_ACCT_STATUS foreign key (STATUS) references APPUSER.ACCOUNT_STATUS (STATUS),
    constraint FK_ACCT_TYPE foreign key (ACCOUNT_TYPE) references APPUSER.ACCOUNT_TYPE (TYPE),
    constraint FK_CUSTOMER foreign key (CUST_ID) references APPUSER.CUSTOMERS (CUST_ID))
    distribute by reference (FK_CUSTOMER);
```

# Using Oracle cx_Python with TimesTen Scaleout



Python [and **Node.js**, **GoLang**, **Ruby** and **PHP**] uses an **OCI driver**
Use tnsnames or easyconnect to connect

**tnsnames.ora** :

sampledb_1812 =(DESCRIPTION=(CONNECT_DATA = (SERVICE_NAME = **sampledb_1812**)(SERVER = **timesten_direct**)))
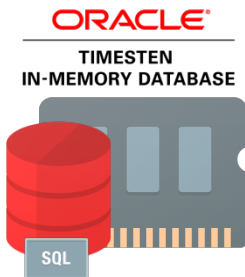sampledbCS_1812 =(DESCRIPTION=(CONNECT_DATA = (SERVICE_NAME = **sampledbCS_1812**)(SERVER = **timesten_client**)))
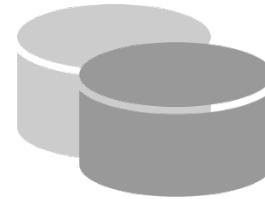
**TimesTen ODBC DSN**

**Client/Server or
Direct Linked**

# Oracle TimesTen In-Memory Database

## Relational Database

- Pure in-memory
- ACID compliant
- Standard SQL
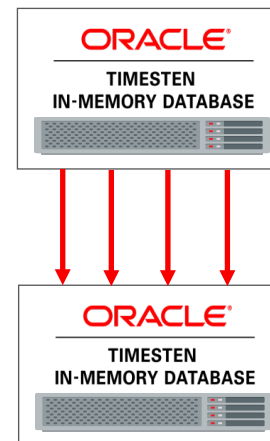- Entire database in DRAM

## Persistent and Recoverable

- Database and Transaction logs persisted on local disk or flash storage
- Replication to standby and DR systems

## Extremely Fast

- Microseconds response time
- Very high throughput

## Highly Available

- Active-Standby and multi-master replication
- Very high performance parallel replication
- HA and Disaster Recovery

**ORACLE**

# Most Widely Used Relational In-Memory Database

**Deployed by Thousands of Companies**

# The Forrester Wave™: In-Memory Databases, Q1 2017

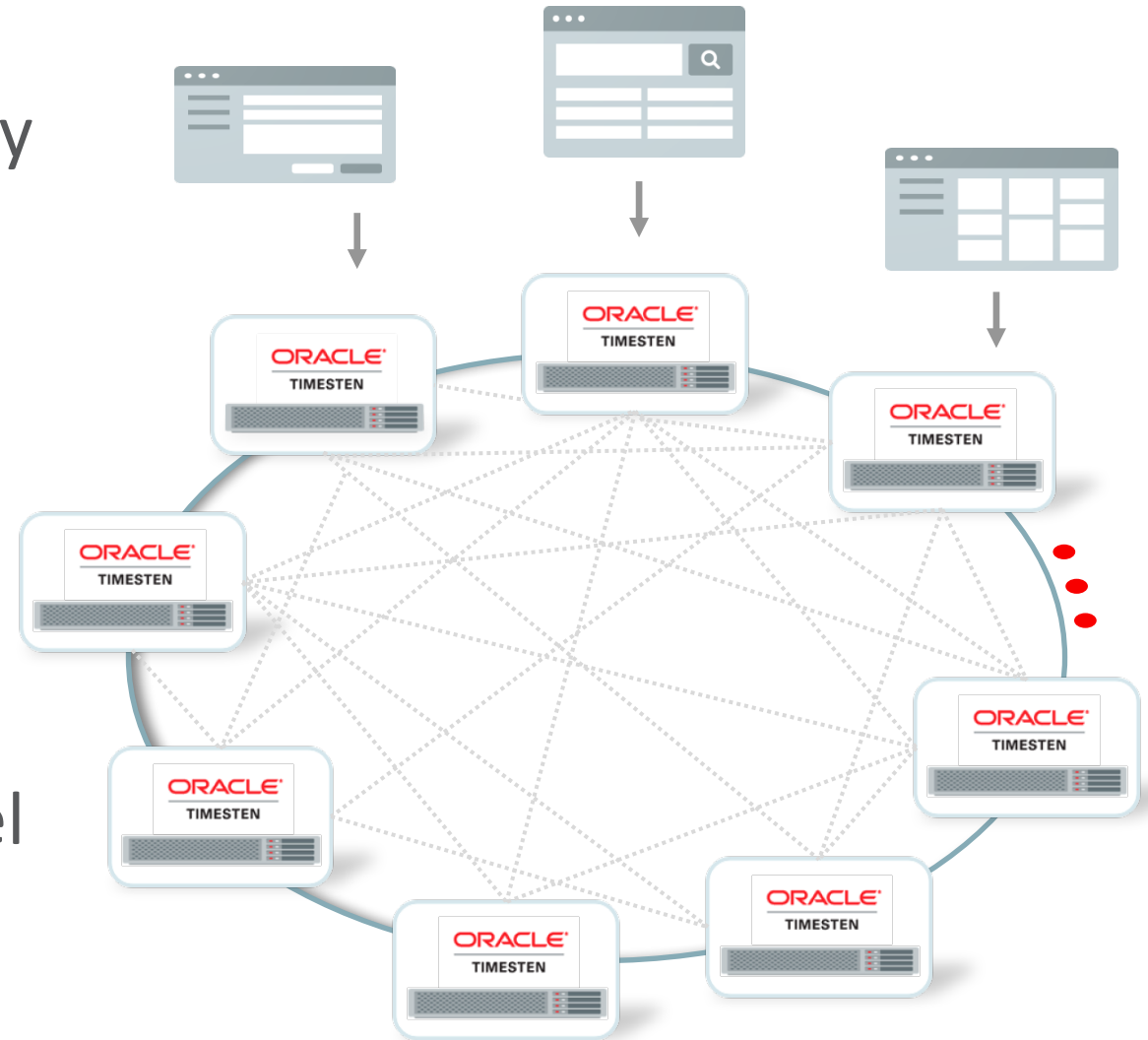**Oracle In-Memory Databases Scored Highest by Forrester** on both Current Offering and Strategy

http://www.oracle.com/us/corporate/analystreports/forrester-imdb-wave-2017-3616348.pdf
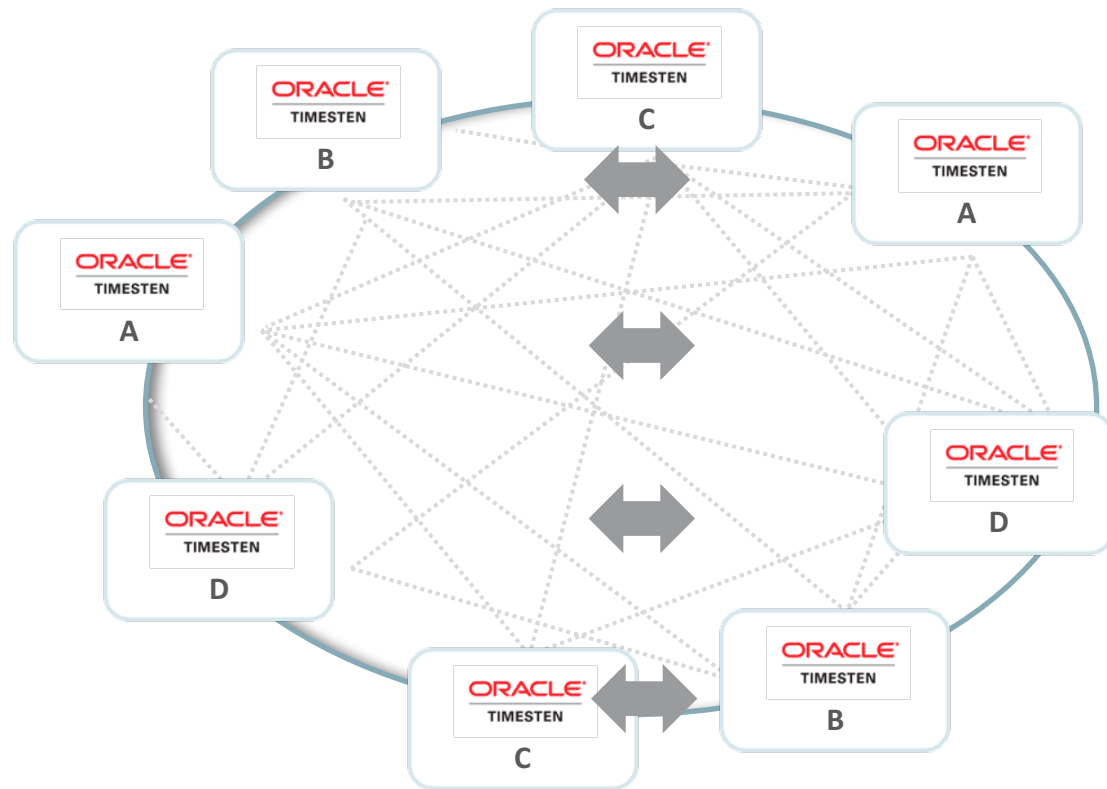
# Single Database Image

- Database size not limited by memory
- Table data distributed across all elements
  - All elements are equal
- Connect to **any** element and access **all** data
  - Distributed queries, joins & transactions
- No need to de-normalize data model

**ORACLE**®

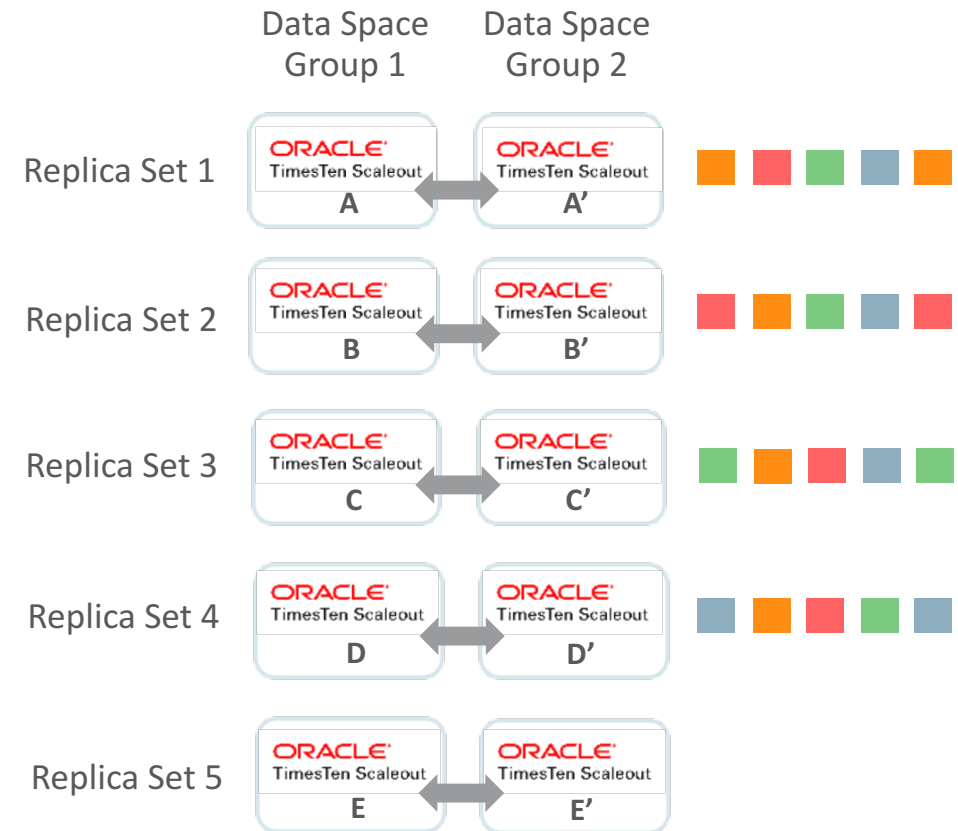# High Availability and Maximum Throughput
**K-Safety, All Active**



- Built-in HA via multiple copies of the data (K-safety)
  - Automatically kept in sync
- **All** replicas are **active** for **reads** and **writes**
  - Double the compute capacity
- Transactions can be initiated from and executed on any replica

# TimesTen Scaleout - Elastic Scalability

Expand and shrink the database based on business needs
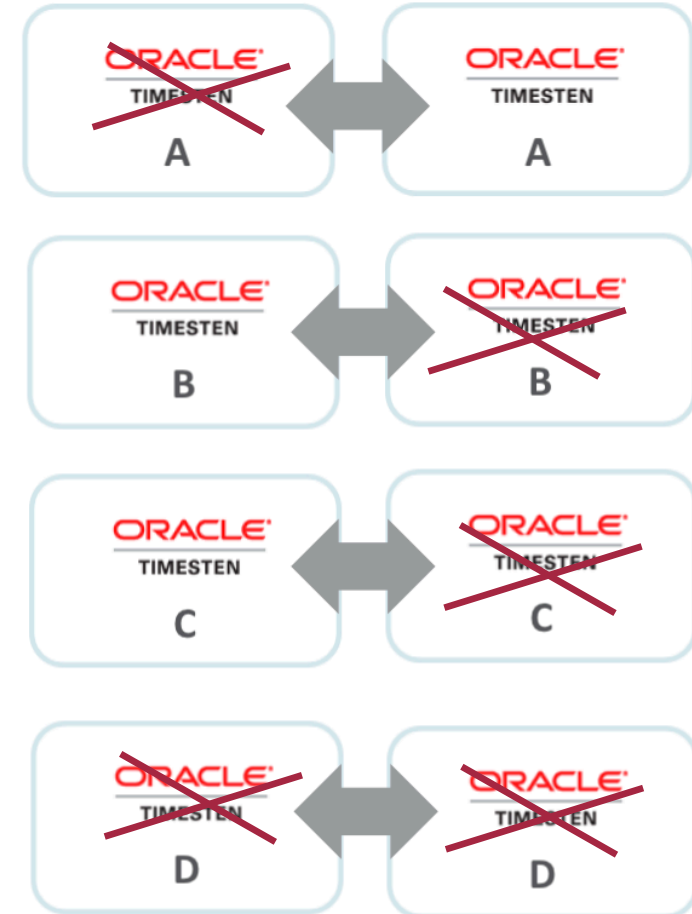
Adding (and removing) database elements

- Data redistributed to new elements
- Workload automatically uses the new elements
- Connections will start to use new elements
- Throughput increases due to increased compute resources

# Database Fault Tolerance – No Application Down Time

**Provided one entire copy of the database is available**

- If multiple elements fail, applications will continue provided there is one complete copy of the database

- Recovery after failure is automatic
- If an entire replica set is down, application can **explicitly** choose to accept partial results