

Oracle Database Reliability

In the context of distributed databases

Francisco Munoz Alvarez

Distinguished Product Manager

Oracle Database High Availability (HA), Scalability and Maximum Availability Architecture (MAA) Team



@fcomunoz



http://www.linkedin.com/in/franciscomunozalvarez



www.oraclemaa.com

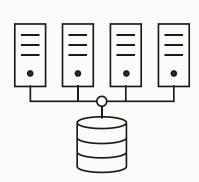


Some key concepts

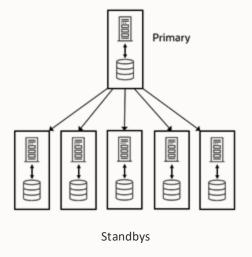


High Availability (HA) and Scalability Concepts

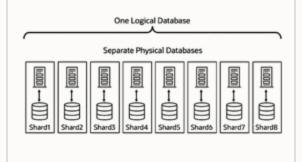
Shared Disk / Shared Cache



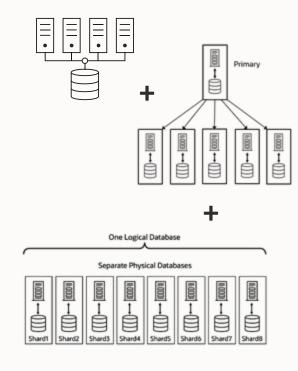
Replication / Read Replicas



Sharding

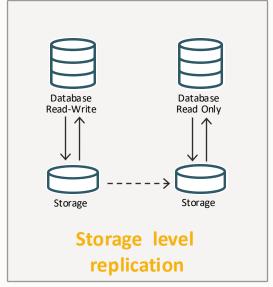


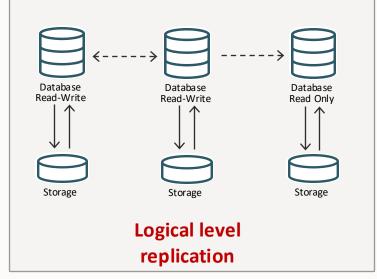
Combination thereof

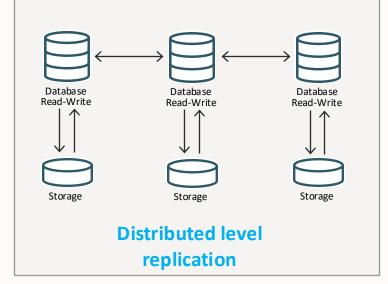




Replication Types







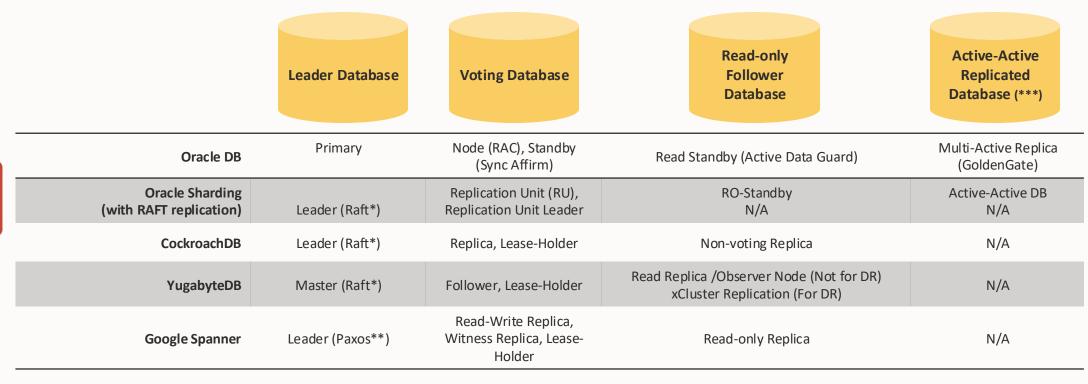
AWS Aurora AlloyDB

Oracle
PostgreSQL
SQL Server

Spanner CockroachDB YugabyteDB TiDB



Same Concepts – Different Names



(***) Asynchronous active-active replication between independent Clusters



New In

23^{ai}

RAFT Protocol in SQL Distributed Databases

RAFT is a consensus algorithm that ensures a group of nodes agree on a series of log entries (like SQL transactions). Each node in a cluster has one of three roles:

- Leader: Handles all client writes.
- **Follower**: Replicates entries from the leader.
- Candidate: Can be elected leader in a new term if the current leader fails.

Key Phases:

- Log replication: Leader appends entries and replicates to followers.
- Leader election: Triggered if the leader is unresponsive. A majority is required for an election.



RAFT Protocol in SQL Distributed Databases

RAFT Leader Election Process:

- 1. A follower detects the leader's unresponsiveness and becomes a candidate.
- 2. It increments the term and sends vote requests.
- 3. Other nodes respond with votes.
- 4. If a majority is reached, the candidate becomes the new **leader**.
- 5. The leader resumes processing log entries.

Impact on Availability:

- During the election, writers are unavailable.
- Short timeouts cause instability; long timeouts delay recovery.



Some SQL Distributed Databases Comparison

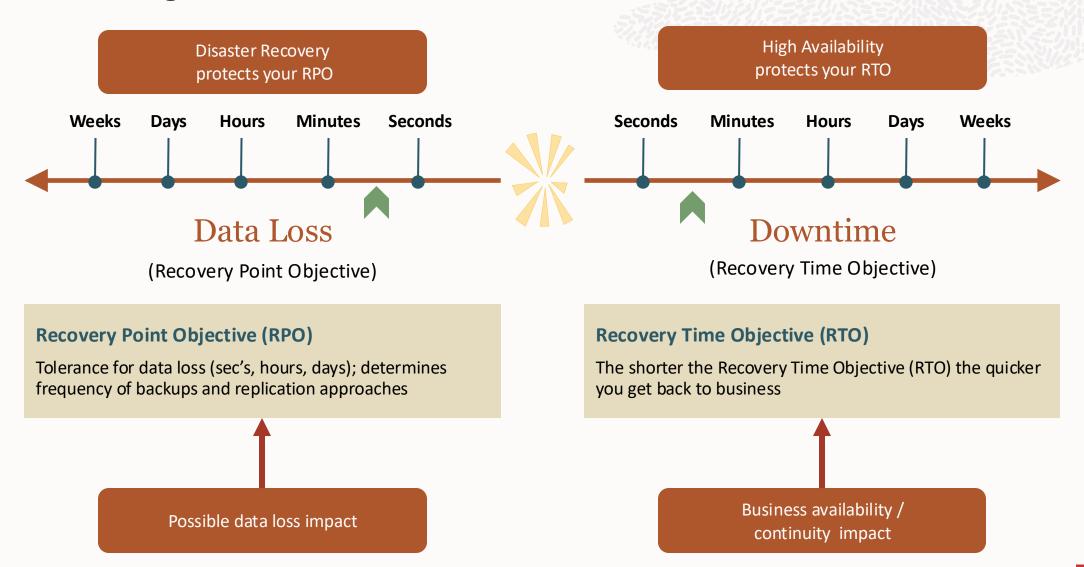
Feature	CockroachDB	YugabyteDB	TiDB	
RAFT Implementation	Custom per-range RAFT	DocDB RAFT with strong sync	Per-region RAFT groups	
Leader Election Time	~300ms (tunable) – whole process up to 9s	~150- 300ms (adaptive) – whole process up to 3s (TCP Timeouts can take up to 15s)	~200ms – whole process up to 20s	
Multi-leader Write Support	No (single per range)	No	No	

Factors Affecting Election Time:

- **Network Latency:** Higher network latency between data centers can increase the time required for leader election and data replication.
- **Cluster Size:** As the number of nodes increases, the overhead of managing elections and replicating logs grows, potentially impacting scalability and election time.



Understanding RPO and RTO



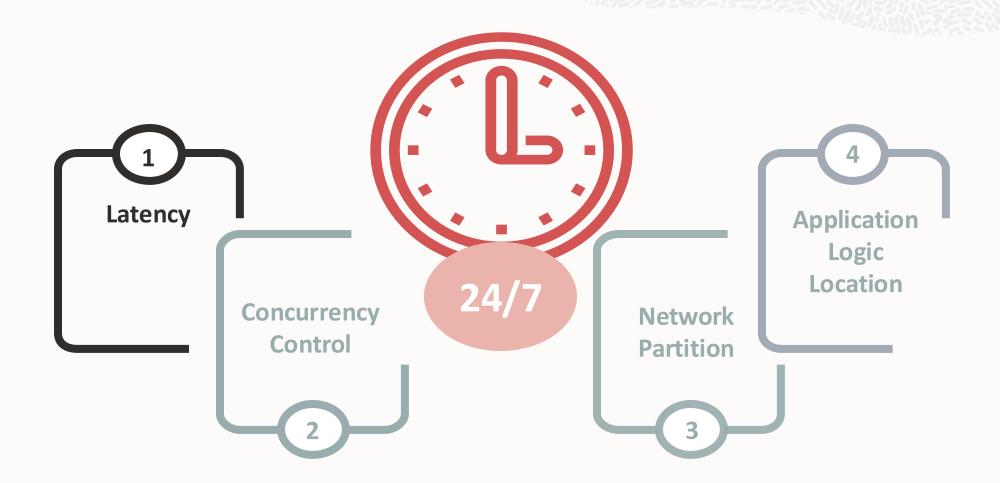
Some More SQL Distributed Databases Comparison

Capability	CockroachDB	YugabyteDB	TiDB	
Read Replica Support	Yes (geo-partitioned)	Yes (follower reads)	Yes (TiFlash replicas)	
Replica Consistency	Eventually consistent	IIMAIINA I ANCICTANAV	Strong Consistency Stale Read (Timestamp Consistency)	
DR Read Usage	Recommended	Recommended	Recommended	

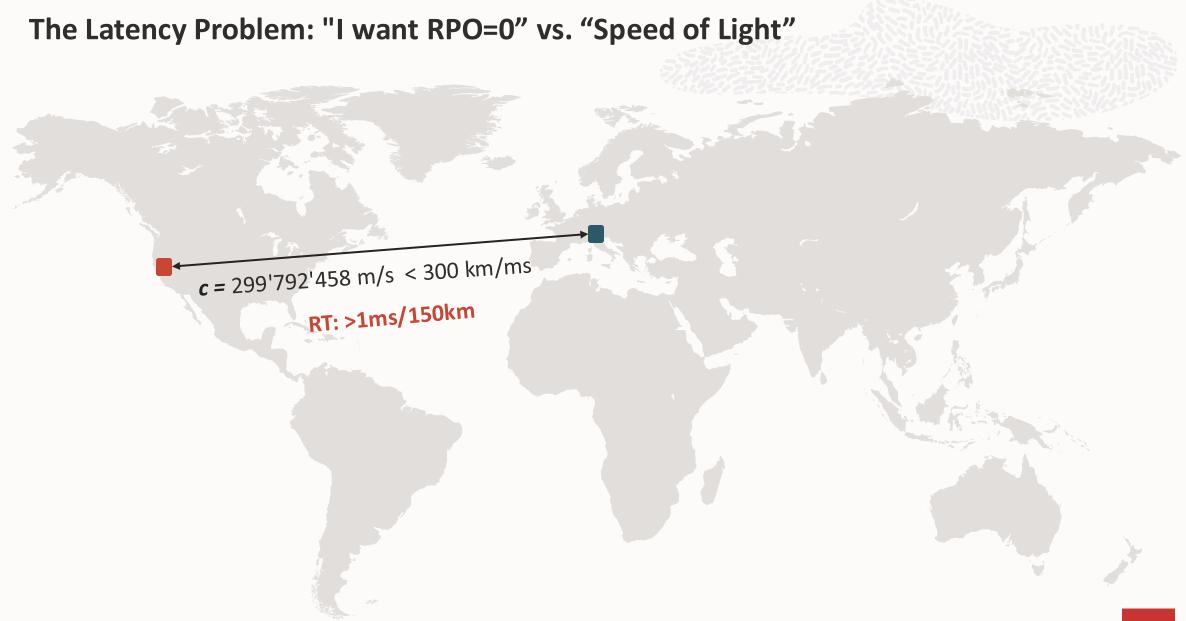


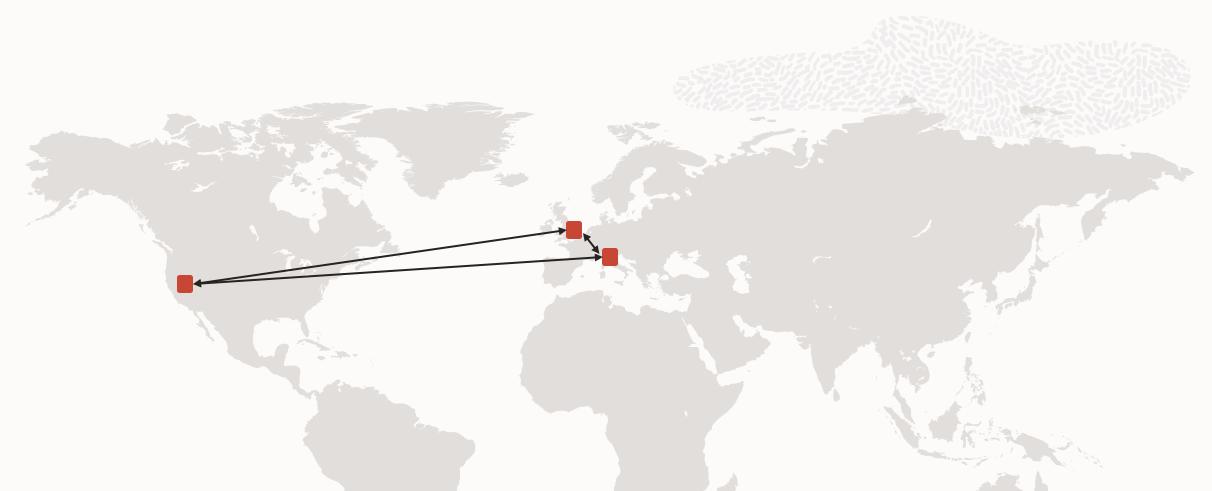
Factors Influencing Data Reliability(*)

Top 4 Factors Influencing Reliability









No distributed database can provide magical capabilities to defeat the speed of light or the harsh reality of cross-region networking at scale.



SQL Distributed Databases Cross-Region Integrated DR/HA

Designing RAFT groups across multiple regions is essential for resilience.

Challenge:

If quorum spans East + West, East outage = no quorum = no writes.

Recommendation:

Split voting nodes across 3 zones and use observer/read-only replicas in DR sites.



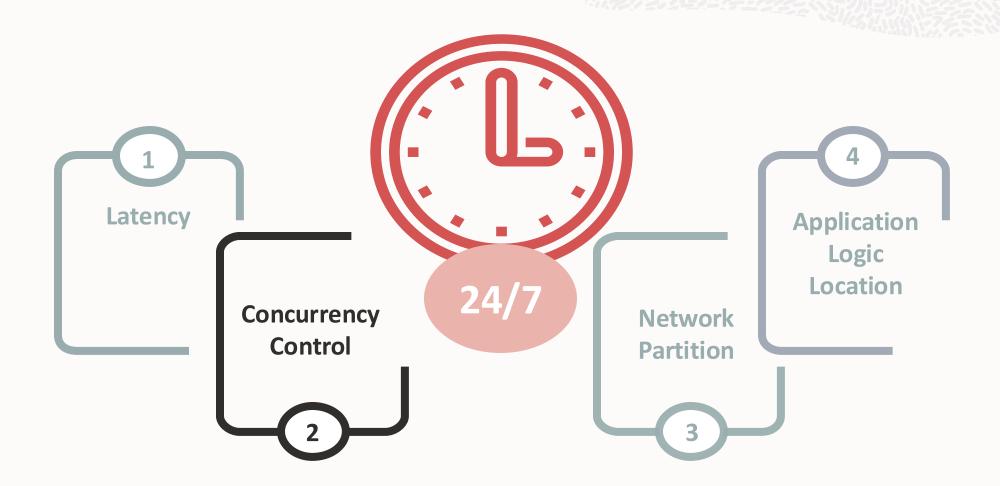
Different Approaches – Different Results

CockroachDB ensures Disaster Recovery (DR) capabilities only within the same region. No multiregion implementations are possible without incurring severe penalties to the applications using it (as per example, latency).

Oracle Database can easily fulfil any cross-region Disaster Recovery (DR) requirements, including multiactive (since 2009) and zero data loss at any distance.



Top 4 Factors Influencing Reliability





Concurrency Control

- Optimistic locking, where a record is locked only when changes are committed to the database
- Pessimistic locking, where a record is locked while it is edited



SQL Distributed Databases Data Consistency & Failover

During failover, transactions may roll back if not committed by a majority.

Feature	CockroachDB	YugabyteDB	TiDB	
Vrite Durability Model Majority write		Majority write	Raft commit based	
Transaction Model (Default)	Serializable	Snapshot Isolation	Snapshot Isolation	



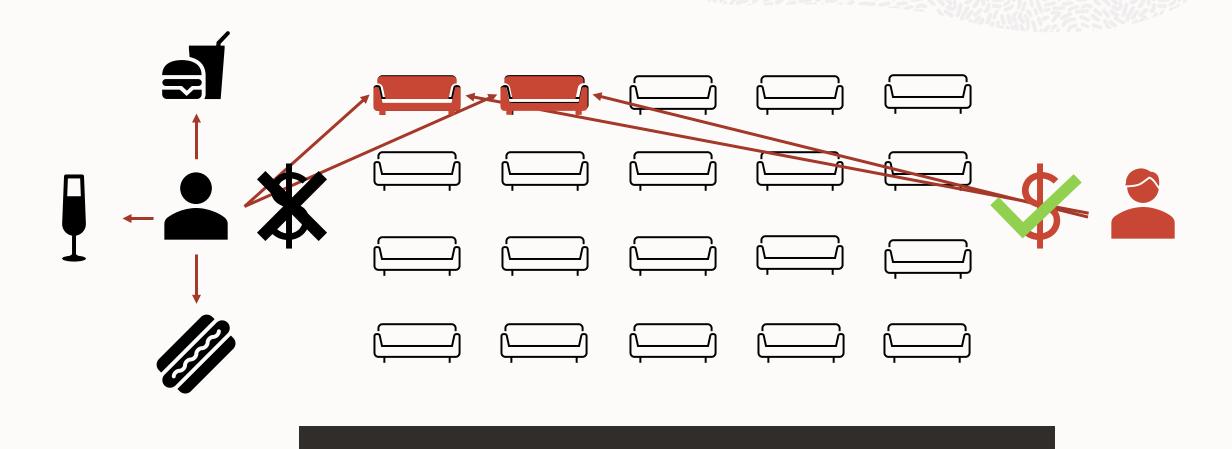
SQL Distributed Databases Concurrency Control and Lost Writes

Most use MVCC and transactional layers on top of RAFT.

Feature	CockroachDB	YugabyteDB	TiDB	
Concurrency Control	MVCC	MVCC + 2PC (two-phase commit)	MVCC + Per-region	
Default Locking Mechanism		Optimistic (Default with Snapshot Isolation level) – Pessimistic (Read Committed Isolation Level)	Pessimistic (Default) and Optimistic	
Lost Write Handling	Strong prevention	Strong prevention	Depends on network	



Example of a hypothetical system using Optimistic Concurrency Control





Different Approaches – Different Results

Using optimistic concurrency control by default (like CockroachDB does) means developers may be required to fully redesign it data model and application

The best of both worlds: Oracle automatically assures read consistency per your requirements (statement-level or transaction-level read consistency).



New with 23ai

Oracle Database SAGA Pattern is a robust approach to maintaining data consistency across
multiple microservices involved in long-running transactions. A saga is a sequence of local
transactions that, collectively, achieve a global transaction. It uses either orchestration or
choreography to manage the sequence and potential rollbacks

Travel Agency App Begin SAGA 1). Book flight If failed Cancel 2) Book hotel If failed Cancel 3) Book car If failed Cancel 4) If all 3 succeed Commit Saga [Id]; Else Abort Saga [Id] (Compensation) End SAGA



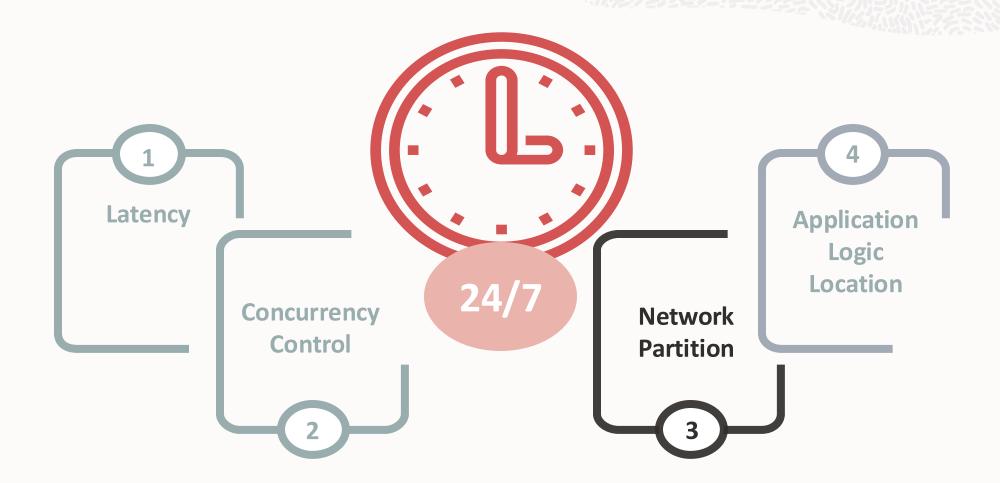
New with 23ai



- Lock-Free Reservation is a feature that significantly improves concurrency on numerical column updates, particularly for "hot" data such as account balances, inventory quantities, and similar numeric aggregate data. It allows multiple concurrent transactions to modify the same column value without being blocked by each other, which is especially useful for scenarios with high concurrency and updates to low cardinality tables.
- Transaction Priority is a feature that allows applications to assign priority levels or levels of importance (HIGH, MEDIUM, or LOW) to transactions. This enables the database to automatically roll back lower-priority transactions that are blocking higher-priority ones, potentially improving overall system concurrency and responsiveness.

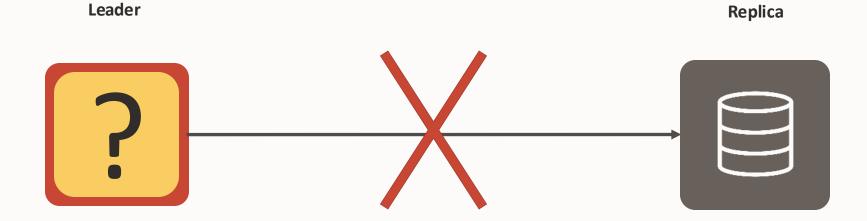


Top 4 Factors Influencing Reliability





Network Partitioning



DID IT CRASH?

DID IT STALL?

DID IT KEEP COMMITTING?

WHAT IS THE LEADER DOING?

- If the leader is still committing
 - DATA LOSS and possible split-brain
- If the leader crashed
 - *Maybe* no data loss



Different Approaches – Different Results

With CockroachDB a network issue and a change of leader could trigger a lost update scenario.

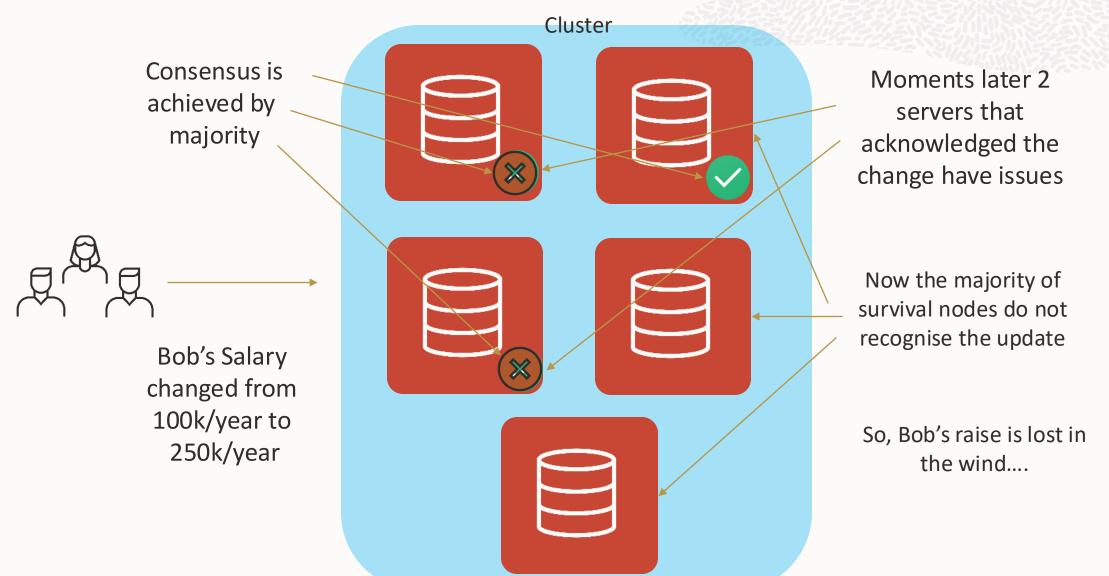
With Oracle Database, lost updates would not occur due to network partitioning.

With Oracle Real Application Clusters (RAC), when used for high availability and scalability

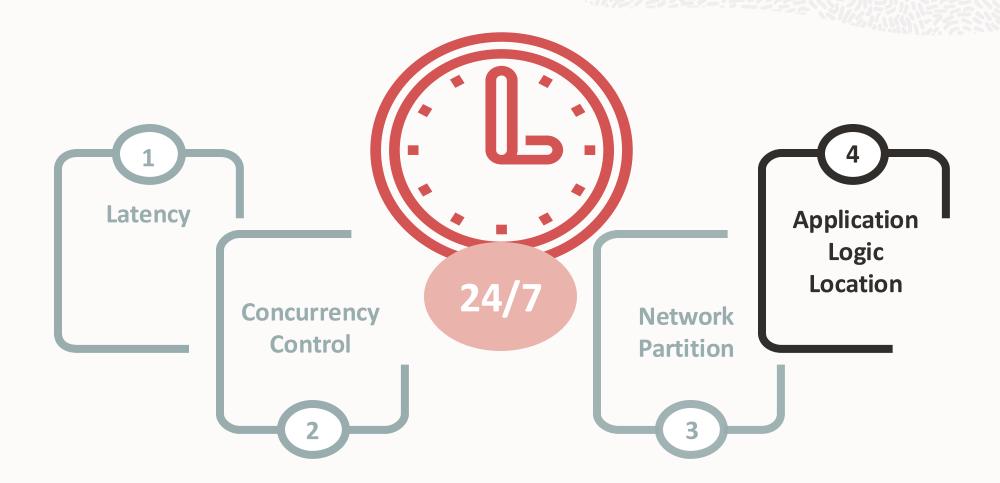
Oracle Data Guard, when used in maximum availability mode or higher



The Hypothetical Journey of a lost update that could happen with Raft replication and optimistic concurrency control



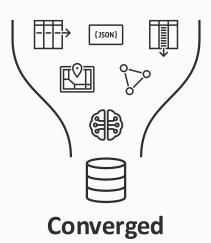
Top 4 Factors Influencing Reliability





Application Logic is best stored in Oracle Converged Database

A complete database that makes it dramatically easier to develop and run modern apps



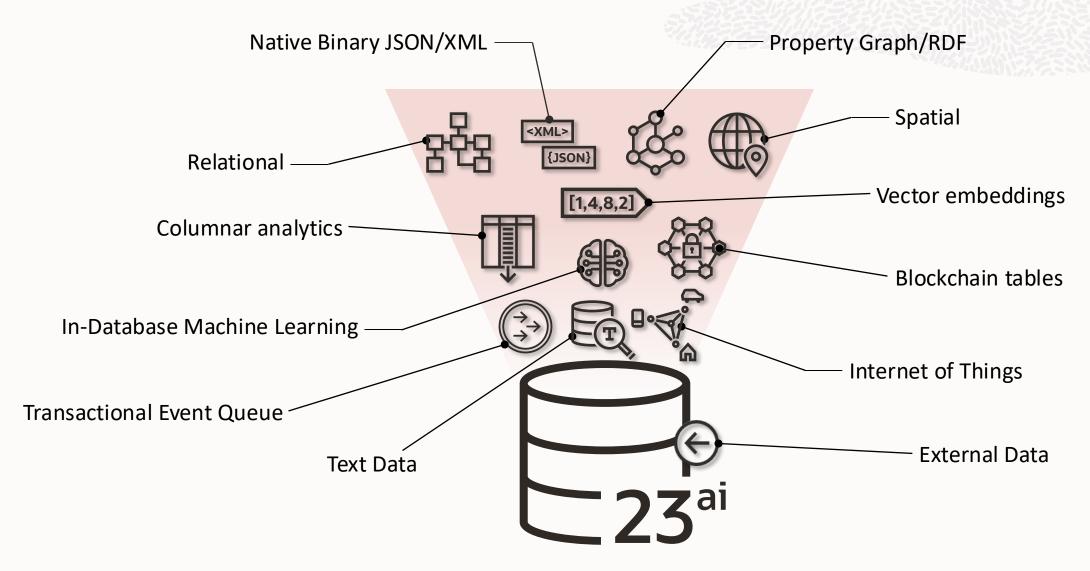
Database

- No need to fragment data across databases to support new app requirements
- Scaling and availability are transparent, without sacrificing data consistency
- No need to compromise on functionality or performance
 - Oracle's data technologies are rated industry-leading in each area

Creating a Fully Complete Database has Taken Decades of Effort by Thousands of Engineers

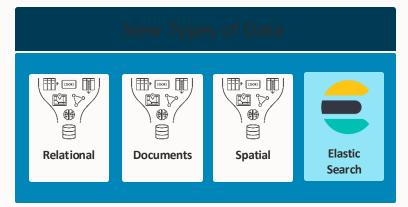


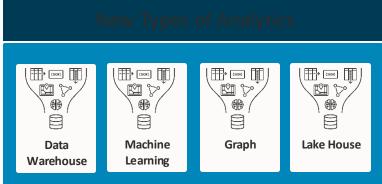
Oracle Converged Database

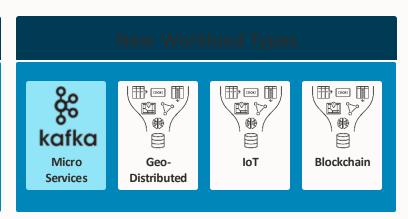


Oracle Converged Database provides choices

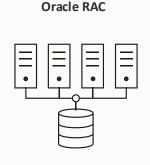
Converged does not mean data must be in one monolith database – YOU choose:

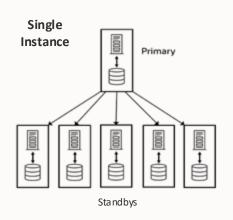


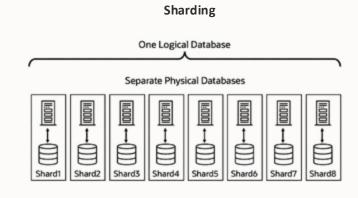




(Support for all modern data types and analytics are included at no additional charge.)









Distributed Databases - Comparison Matrix

	On-Premises Availability	Distribution Layer	Data Modeling	Converged Database	Triggers, Cursors and Stored Procedures within DB	Distributed Capabilities Since	CC (*) Mode (Default)	Popularity Ranking**
CockroachDB	Yes	Raft with Range Sharding (Hash- Sharded Indexes)	Limited Data and Partition Types	No	No	2015	Optimistic	64
YugabyteDB	Yes	Raft with Hash and Range Sharding	Limited Data and Partition Types	No	Yes	2016	Optimistic	111
Google Spanner	No	Raft with Range Sharding	Limited Data Types	No	No	2017	Pessimistic	97
Oracle Database	Yes	In-Memory Redo (Physical and Logical) Replication, Shard (*), and Raft NEW IN 23ai	Huge amount of Data and Partition Types	Yes		1986 with version 5 Many more years an all other comp		1

^(*) Concurrence Control ** Reference: https://db-engines.com/en/ranking



^{(*) &}lt;u>User-Defined</u>, <u>Range-Based</u>, <u>List-Based</u>, and <u>Composite Sharding</u>.

Oracle Database with MAA



Impact of downtime



\$350K

average cost of downtime per hour



87 hours

average amount of downtime per year



\$10M

average cost of unplanned data center outage or disaster



91%

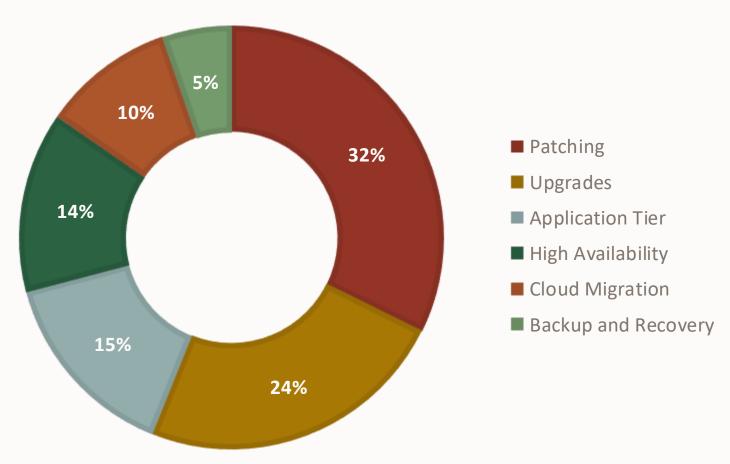
percentage of companies that have experienced an unplanned data center outage in the last 24 months



Addressing Day to Day Challenges...

- Minimize downtime and complexity during patching, migrations & upgrades
- Eliminate data loss and downtime during unexpected outages
- Application tier connectivity and role transitions
- Be prepared while avoiding complexity

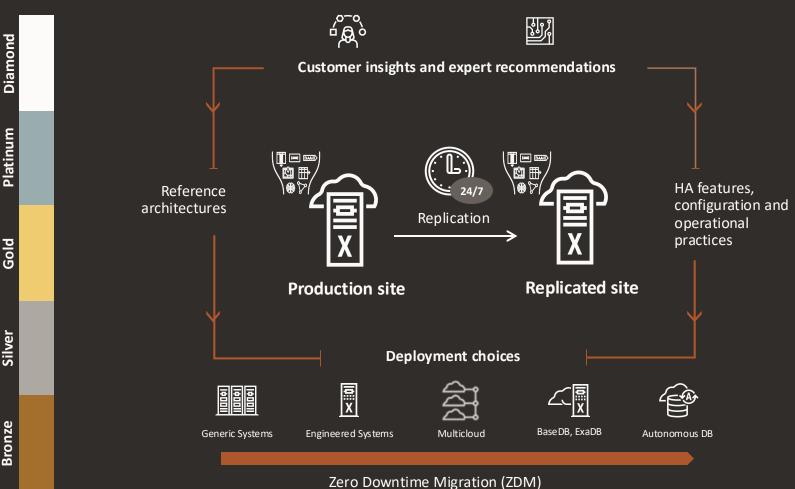
TOP PAIN POINTS*







Next Gen Maximum Availability Architecture (MAA)





Application Continuity



Data protection







ZDLRA+ ZRCV

Active replication







Active Data Guard

Scale out & Lifecycle











Database

Application Testing

Next-Gen MAA Reference Architectures

Availability service levels for the next generation of Oracle Database

Bronze Silver Gold Platinum Diamond (NEW)

Dev, test, prod

Single instance DB

Restartable

Backup/restore



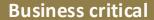
Prod/departmental

Bronze +

Database HA with RAC or Local Data Guard

Client failover HA best practices

Application Continuity (optional)



Silver with RAC+

DB replication with (Active) Data Guard with automatic failover

Client failover DR best practices



Mission critical

Gold with Exadata and either option:

Option 1: GoldenGate with Oracle Database 19c

OR

Option 2: (Active) Data Guard with Oracle AI Database 26ai



Extreme availability

Configuration

GoldenGate 26ai replicas, each running:

Oracle Al Database 26ai

+ RAC on Exadata

+ (Active) Data Guard



Recoverable local failure:

Minutes to hour

Disasters: Hours to days

RPO < 15 min

Recoverable local failure: seconds to minutes
Disasters: Hours to days

RPO < 15 min

Recoverable local failure: Less than 60 seconds Disasters: < 5 min RPO = zero or near zero Recoverable local failure: Less than 20 seconds
Disasters: < 30 secs
RPO = zero or near zero

Recoverable local failures: Less than 10 seconds Disasters zero to 10 secs RPO = zero or near zero

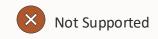


Oracle Database with MAA vs CockroachDB at a Glance



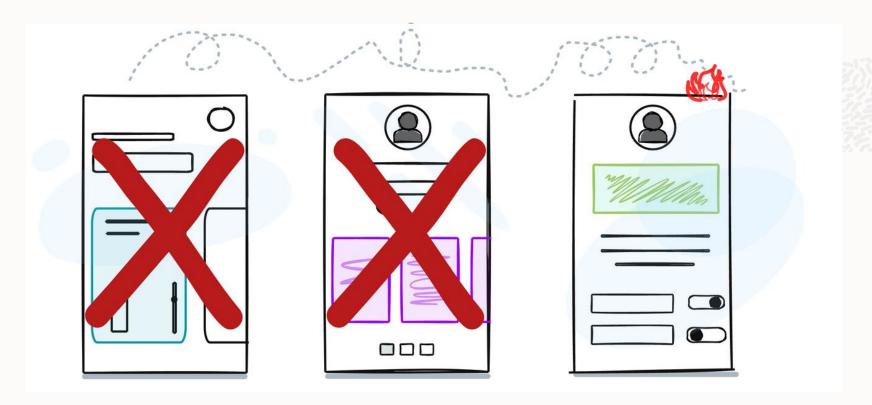
Capabilities	Oracle Gold & Platinum MAA	CockroachDB
Full Disaster Recovery (regional protection)		•
High Consistency (within a region)		
One Node Survival		×
Time Machine Recovery		×
Full Backup and Recovery		lacksquare
Full Backup and Recovery (with Zero Data Loss)		×
Near Zero Downtime Upgrade/Patching		
Automatic Conflict Detection and Resolution		×
Ransomware Protection at Database Level		×
Automated Failover and Transaction Replay		×
End-to-End Validation – Data Corruption Prevention		×
Zero Data Loss at any Distance		×
Application Versioning		×
Converged Database		×
Data sovereignty		lacksquare
Easy to Install/Configure	×	









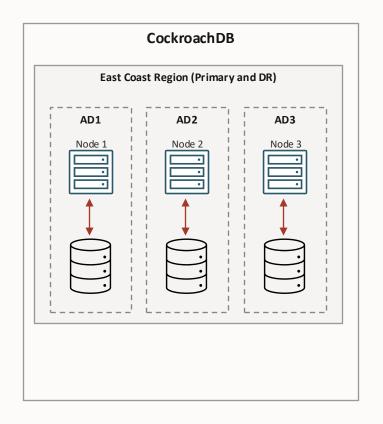


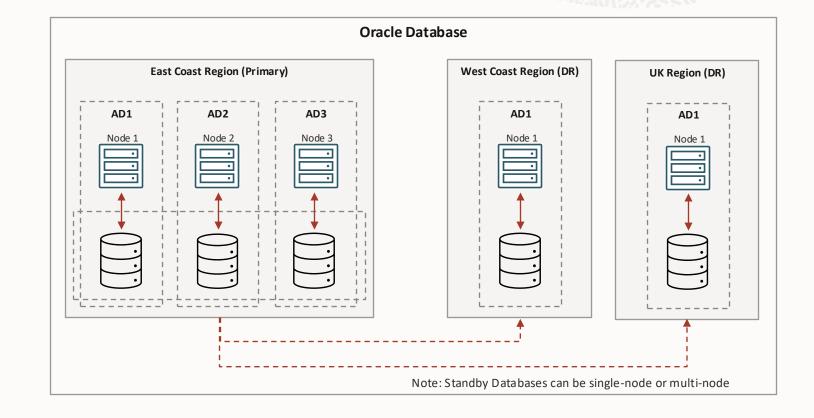
No deliberation and interest of the section of the

A distributed database is a useful technology for achieving an acceptable level of high availability but is *not* built with disaster recovery in mind.



Disaster Recovery (DR) Scenarios







Industry Leading



Leader in Translytical Database.

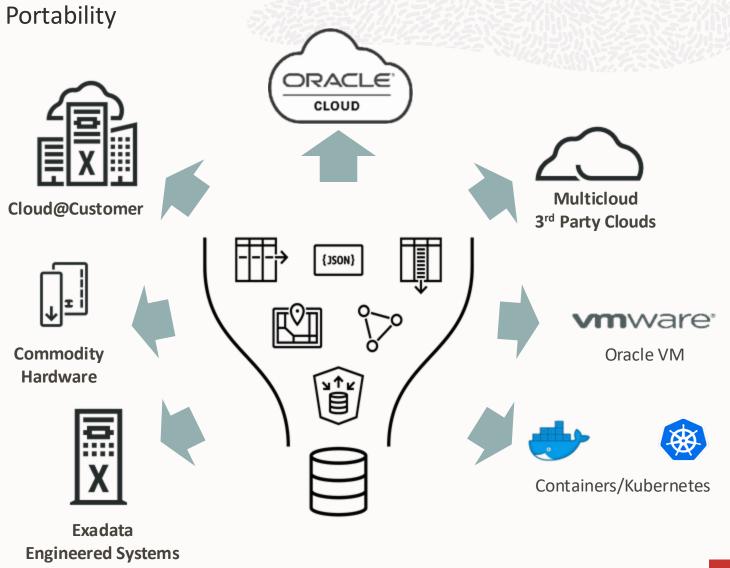
Forrester rates Oracle the strongest leader in the Forrester Wave: Translytical Data Platforms Q4 2024. The rating in this Forrester evaluation validates Oracle Database for its ability to support converged OLTP and Analytical (Translytical) workloads using Database In-Memory, Exadata, extensive multi-model capabilities, and support for relational and non-relational, structured and unstructured.



Oracle Database – Deployment Choices

Deploy Oracle Anywhere – Extreme Portability

Same database, same skills Fast and efficient deployments





Conclusion

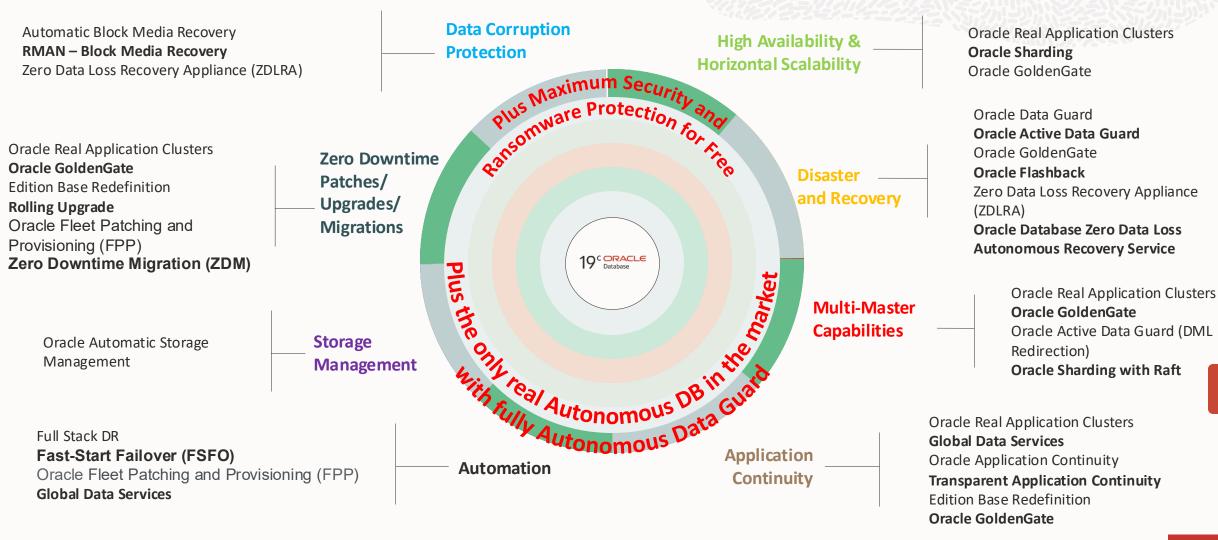


When choosing a technology for distributed database deployments it is important to balance the architecture requirements with the business needs and service level objectives



Summary and Conclusion

Oracle Database is all-inclusive, converged by design, and protected by MAA





NEW IN

23c

Top 5 reasons you should always take into consideration when thinking to replace Oracle with a distributed database

- 1. Changes to the application and data model layer are required to handle requirements that were naturally managed by the database.
- 2. You may need multiple technologies to fulfill all business/application requirements (adding complexity, risk, and cost).
- 3. A "low-cost solution" comes with a cost, including a lack of crucial functionalities and being charged for Enterprise Support and backup, not free.
- 4. The "try again later, scenario" A lost transaction equals lost revenue and, maybe, brand impact!
- 5. Finally, possible data loss due to inefficient Disaster Recovery capabilities. Ensure all your requirements can be fulfilled, especially when discussing RTO and RPO requirements.





Want more details?

- The Fundamentals of Reliability in the Context of Distributed Databases Part
 1
- From Chaos to Order: The Importance of Concurrency Control within the Database Part 2
- <u>The Need for Speed: Physics, Latency, Emergent Technologies, and the Future</u> of Enterprise Database Transactions - Part 3
- Demystifying the Safety Net: Is Your Database Prepared for Anything? Part 4



Any Questions? Thank you!



@fcomunoz



http://www.linkedin.com/in/franciscomunozalvarez



www.oraclemaa.com

