# Studying state-of-the-art HTAP systems

- Satya Sai Bharath Vemula

- Rwitam Bandyopadhyay

- Shubham Pandey

# Table of Contents

**Introduction (~4 min)**

- History of OLAP & OLTP

- Challenges of HTAP Systems
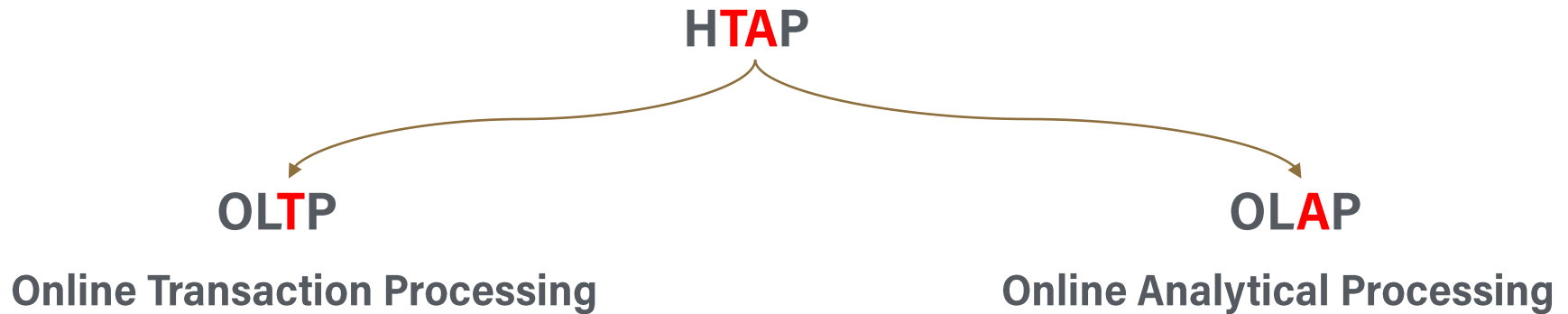
- Architectural taxonomy of present HTAP Systems

**Oracle Dual-Format Database (~8 min)**

**SAP HANA (~8 min)**

**TiDB (~9 min)**

**Project Plan (~1 min)**

# HTAP

## OLTP

**Online Transaction Processing**

## OLAP

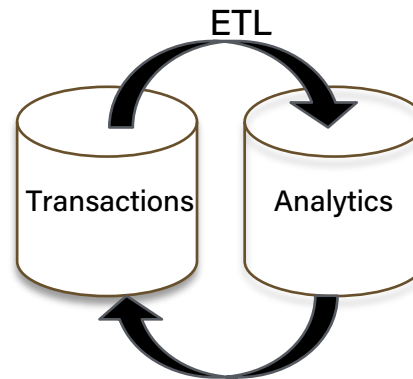**Online Analytical Processing**
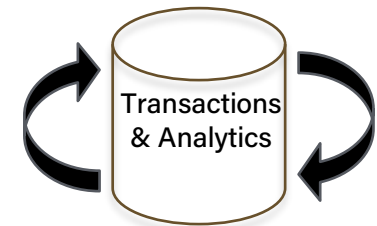
ETL

# Challenges of HTAP Systems

## What is it not easy to imagine both systems in one?

- Data modeling complexity

- Performance optimization
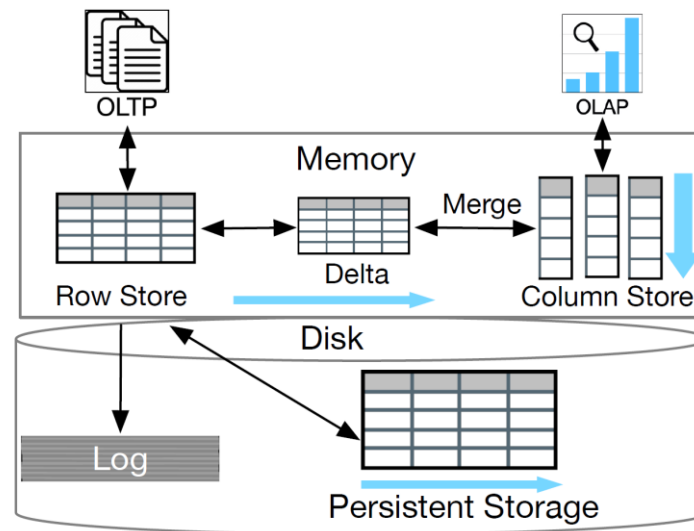
- Consistency and correctness

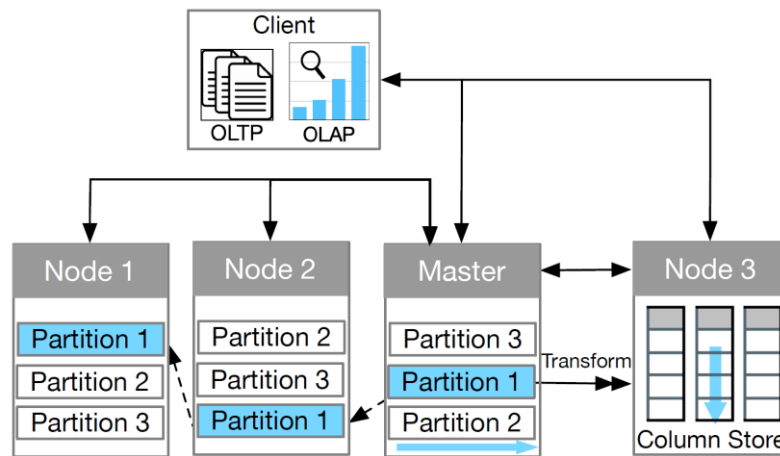- Scalability

- Cost

Traditional Systems

HTAP Systems

ETL

Transactions

Analytics

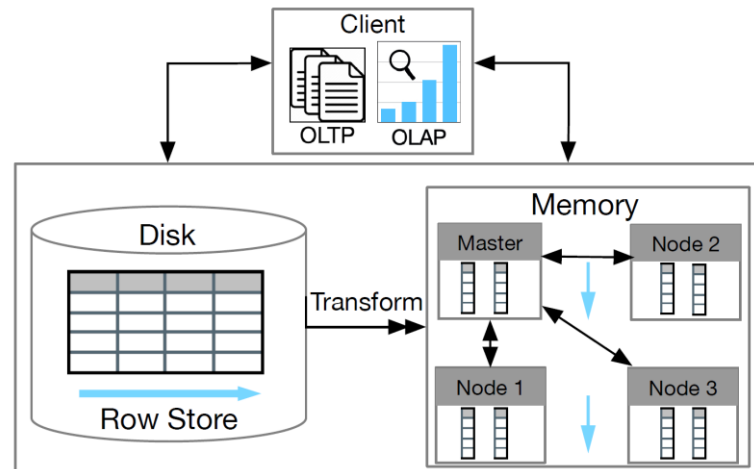Transactions & Analytics

**Primary Row Store + In-Memory Column Store**

Distributed Row Store + Column Store Replica

Source: Guoliang Li and Chao Zhang. 2022. HTAP Databases: What is New and What is Next. (SIGMOD '22)

# Storage Architectures of HTAP Databases



Disk Row Store + Distributed Column Store

Source: Guoliang Li and Chao Zhang. 2022. HTAP Databases: What is New and What is Next. (SIGMOD '22)

**Primary Column Store + Delta Row Store**

# Oracle Dual Format

## Disagg1

Presented by Rwitam Bandyopadhyay
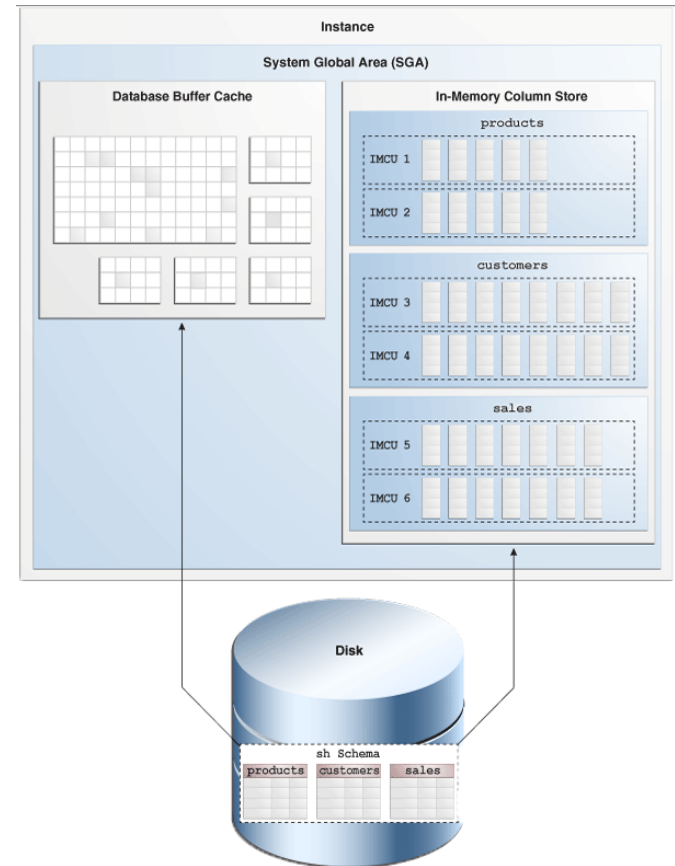
**PURDUE UNIVERSITY®**

# *Motivation*

**Oracle introduced the Database In-Memory option in 2014 as the industry's first dual-format, in-memory RDBMS**

- No single data format is ideal for all type of workloads

- "Primary Row Store + In-Memory Column Store" architectural style

- The new columnar format is a pure in-memory format with no impact to the disk representation

- Tables required for fast analytics can be populated into the In-Memory column store

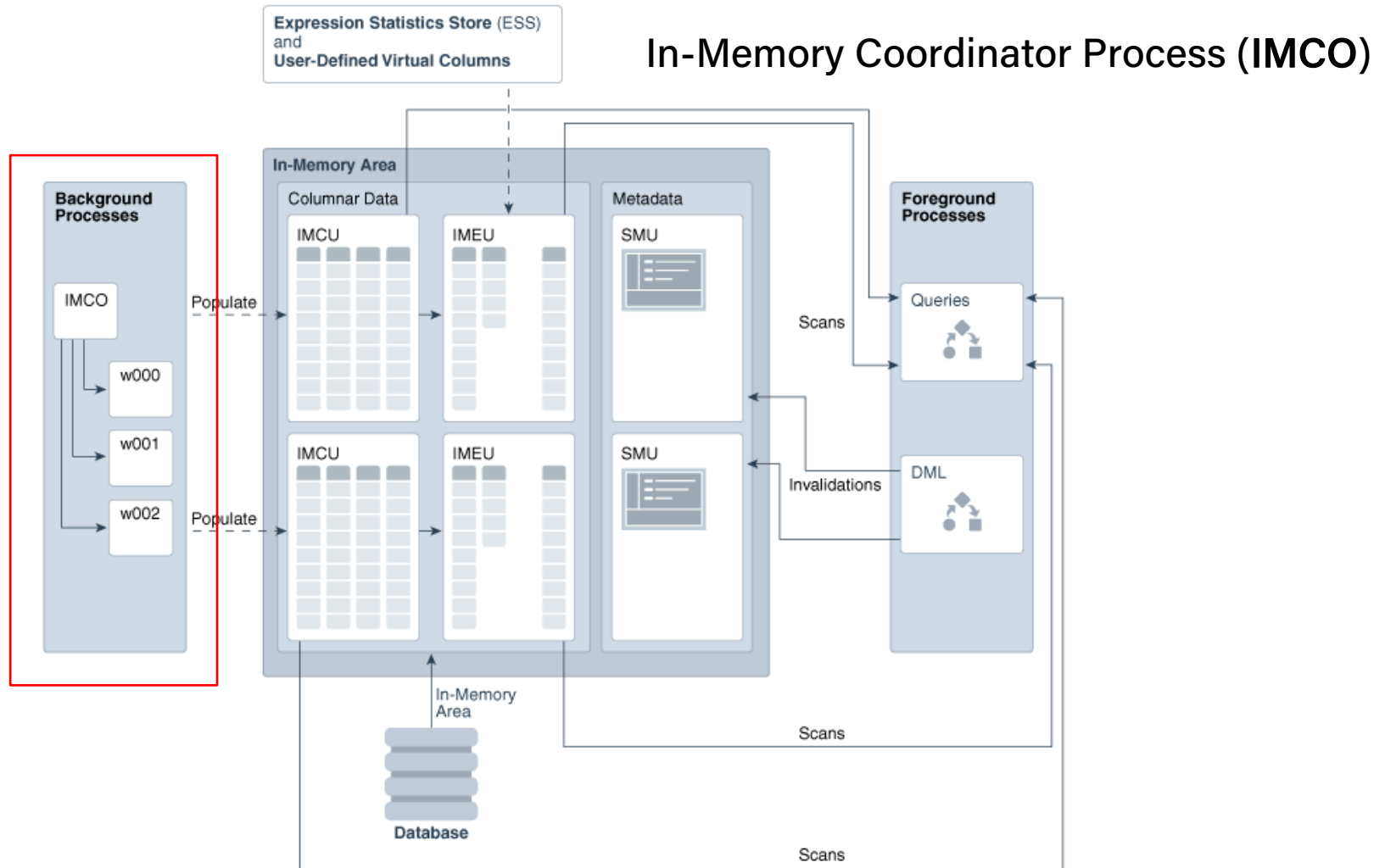- Can be plugged in to any existing system with zero changes

# *Dual-Format, Dual Memory?*

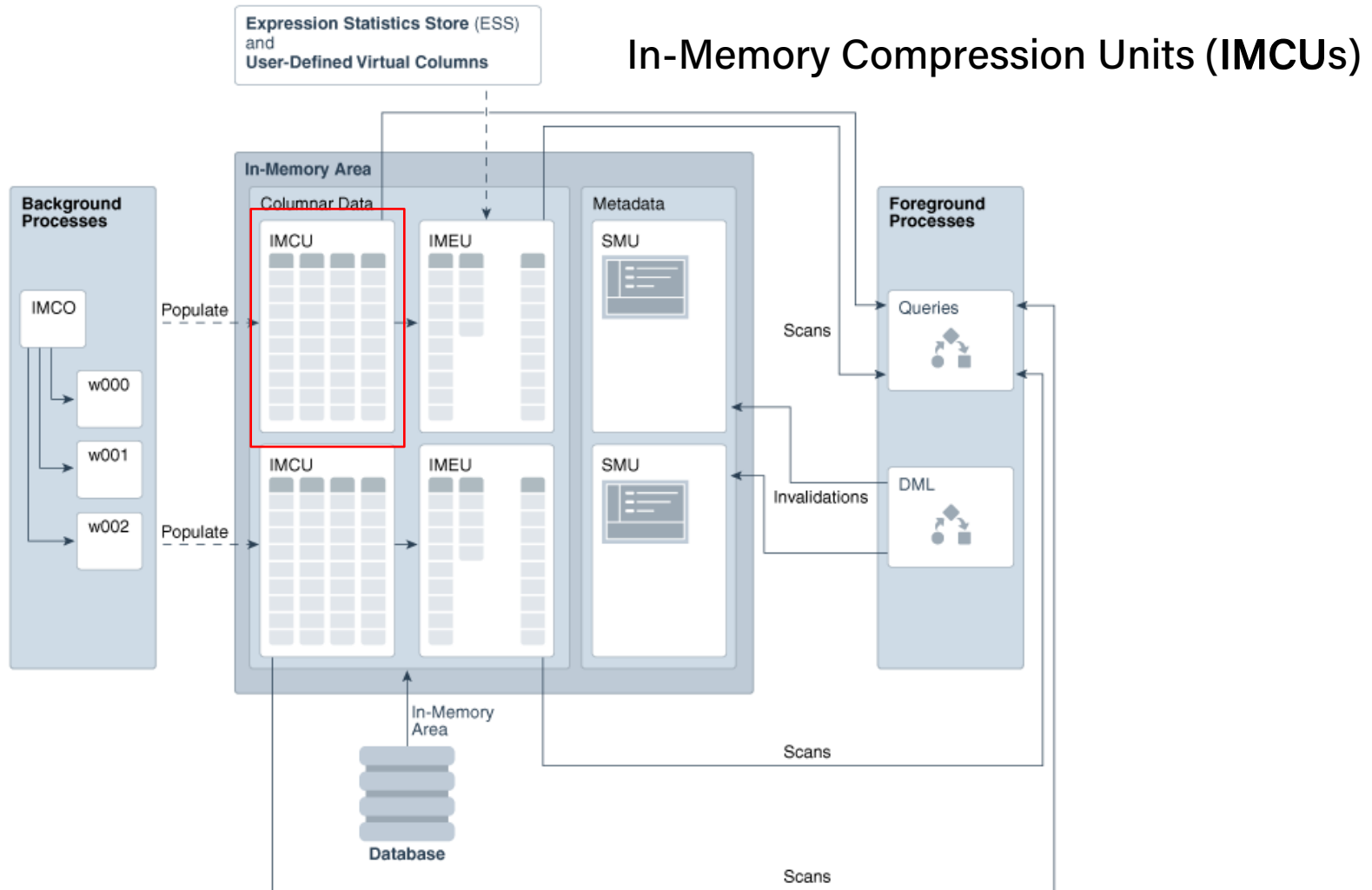**Dual format representation doesn't double memory requirements!**

- Tables required for fast analytics
  → In-Memory Store

- OLTP updates/ highly selective lookups
  → Buffer Cache !

- Analytic Indexes can be dropped from
  the DB

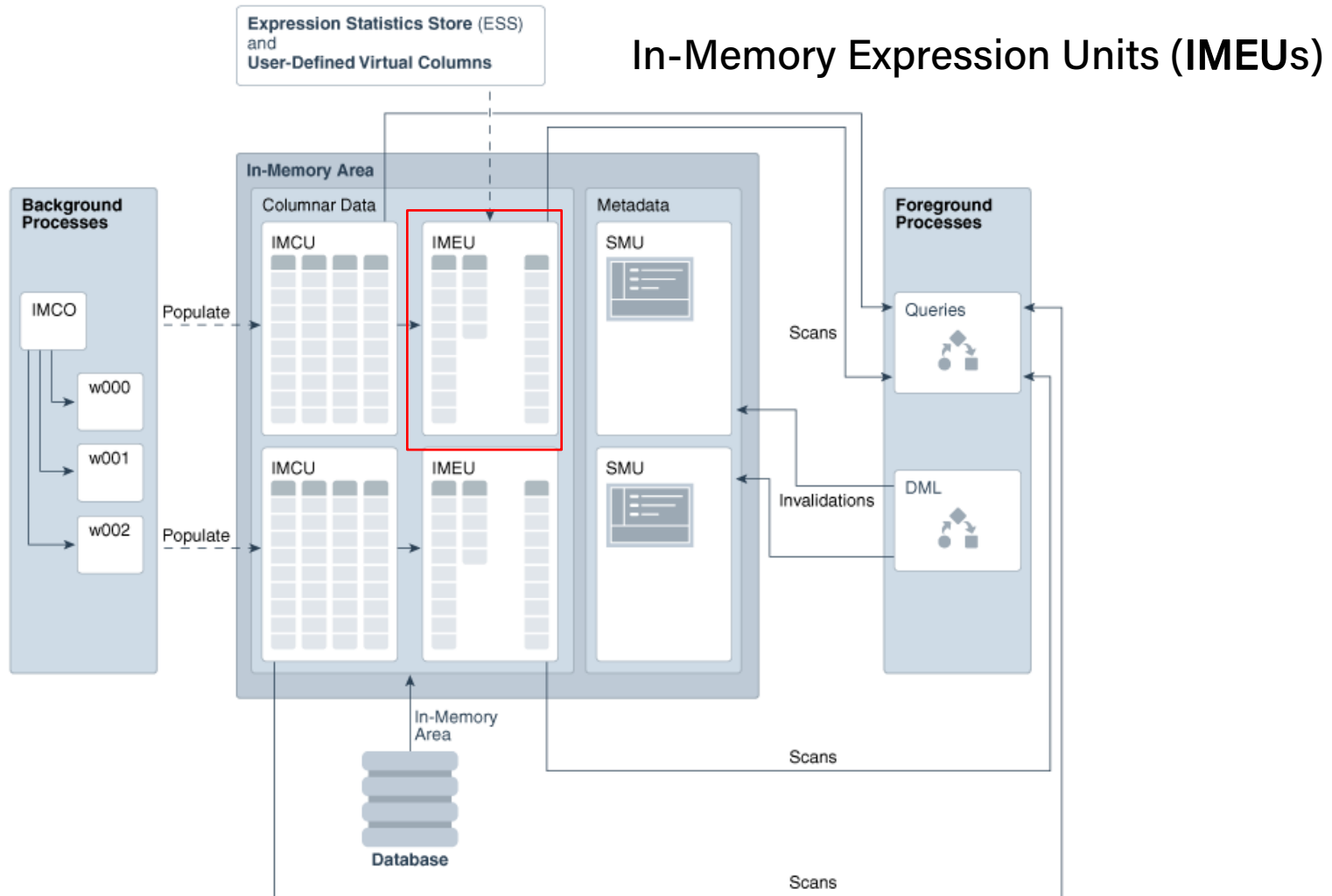- The IM column store is built into the DB
  data access layer



Source: <u>Oracle Database R21 In-Memory Column Store Architecture</u>

# *In-Memory Columnar Format (Building Blocks)*



In-Memory Coordinator Process (**IMCO**)

Source: Oracle Database R21 In-Memory Column Store Architecture

# In-Memory Columnar Format (Building Blocks)

In-Memory Compression Units (**IMCUs**)

# *In-Memory Columnar Format (Building Blocks)*



In-Memory Expression Units (**IMEUs**)

Source: Oracle Database R21 In-Memory Column Store Architecture

# In-Memory Columnar Format (Building Blocks)



Snapshot Metadata Units (**SMU**s)

# In-Memory Columnar Format (Data Population)

**It is possible to selectively populate the IM column store with a chosen subset of the database**

- IM column store is populated using a pool of background server processes

- No application downtime during table population since it continues to be accessible via the buffer cache.

- Most other in-memory databases require a wait time for all objects to be brought in-memory

| Priority | Description |
|---|---|
| CRITICAL | Object is populated **immediately** after the database is opened. |
| HIGH/ MEDIUM/ LOW | Objects are populated after CRITICAL in order of priority. |
| NONE | Object is only populated after it is **scanned for the first time** (Default). |

Source: T. Lahiri, et al. Oracle Database In-Memory: A Dual Format In-Memory Database. In 2015 IEEE 31st International Conference on Data Engineering

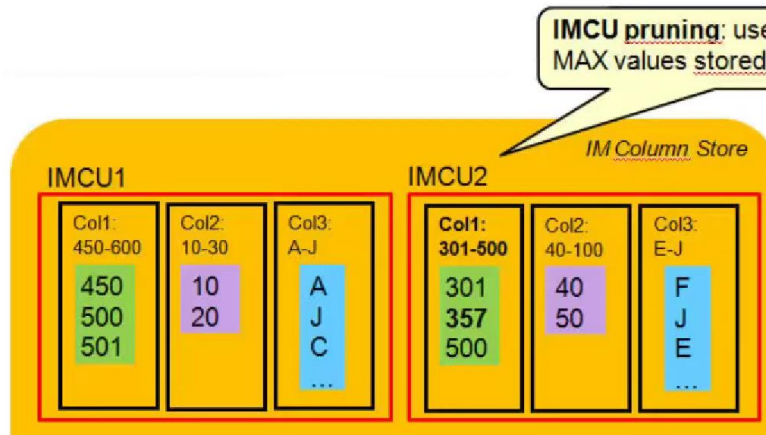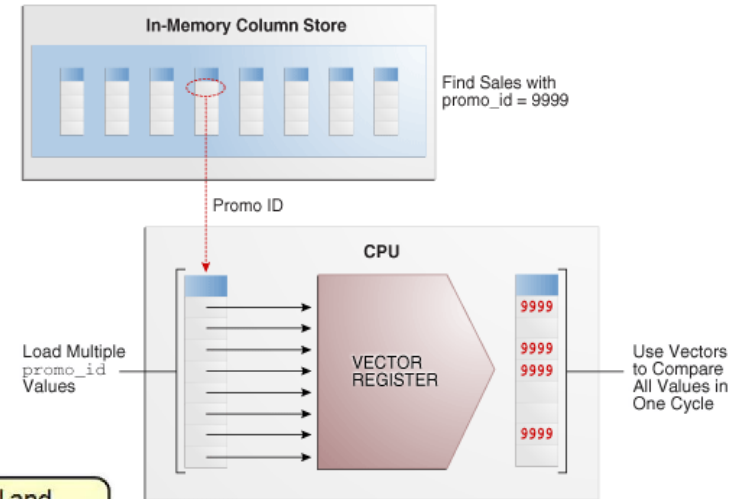# In-Memory Columnar Format (Data Compression)

**It is possible to selectively populate the IM column store with a chosen subset of the database**

- In-memory storage is populated in contiguously allocated units called **In-Memory Compression Units.**

- It is possible to use different compression levels for different partitions within the same table. For example, a SALES table might have →
  - Current week partition - DML Compression
  - Earlier year partitions - Query Compression
  - Decade old partitions - Capacity Compression

| MEMCOMPRESS | Description |
|---|---|
| DML | Minimal compression optimized for DML Performance |
| QUERY LOW/HIGH | Optimized for fastest query performance (Default) |
| CAPACITY LOW/HIGH | Optimized for space saving |

Source: T. Lahiri, et al. Oracle Database In-Memory: A Dual Format In-Memory Database. In 2015 IEEE 31st International Conference on Data Engineering
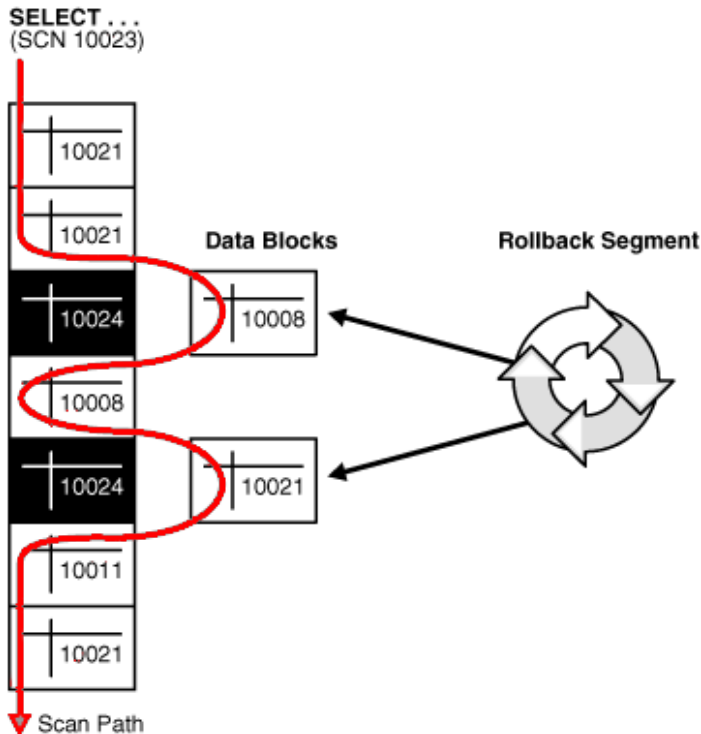
# Query Processing

- In-Memory Scans utilize **SIMD Vector Processing**

- **In-Memory Storage Indexes** are automatically created and maintained on each of the columns.



**IMCU pruning**: use of MIN and MAX values stored in each IMCU



These allow data pruning to occur based on the filter predicates supplied in a SQL statement.

# *Transactional Consistency*

**The default isolation level of Oracle DB is "Consistent Read"**



SELECT . . .
(SCN 10023)

10021
10021
Data Blocks    Rollback Segment
10024    10008
10008
10024    10021
10011
10021
Scan Path

- As a query enters the execution stage, the current system change number (**SCN**) is determined.

- Each query returns all committed data with respect to the SCN recorded at the time that query execution began.
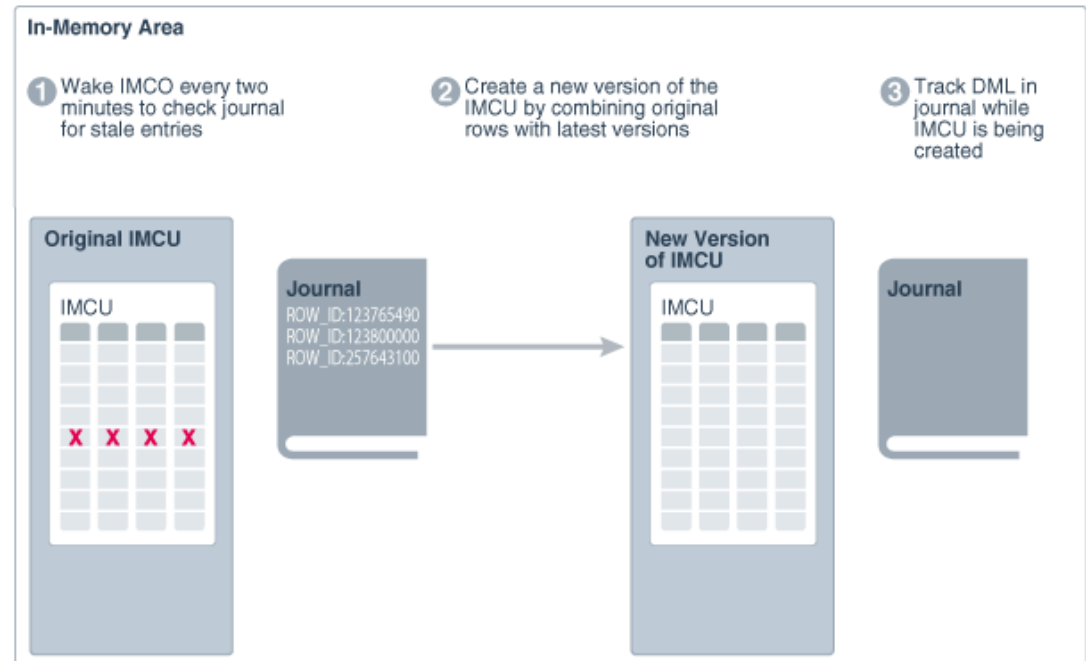
- The In-memory Column store follows similar semantics!

Source: <u>Oracle Database Concepts 10g - Data Concurrency and Consistency</u>

# Transactional Consistency (Double Buffering)

Source: Oracle Database R21 In-Memory Optimizing Repopulation of the IM Column Store

There are two types of strategies for data repopulation:

**In-Memory Area**

① Wake IMCO every two minutes to check journal for stale entries

② Create a new version of the IMCU by combining original rows with latest versions

③ Track DML in journal while IMCU is being created

**Original IMCU**

IMCU

**Journal**
ROW_ID:123765490
ROW_ID:123800000
ROW_ID:257643100

X X X X

**New Version of IMCU**

IMCU

**Journal**

- **Threshold Repopulate**

- **Trickle Repopulate**

At most $\frac{1}{10}^{th}$ of a single CPU core is dedicated to trickle repopulate

Source: Oracle Database R21 In-Memory Optimizing Repopulation of the IM Column Store

# In-Memory Column Store Scale Out



## Distribution

## Duplication

**Distributing Partitions by Hash**

**Distributing Partitions by Range and Sub partitions by Hash**

# In-Memory Column Store Scale Out (Distribution)

**Distribution by Rowid Range**
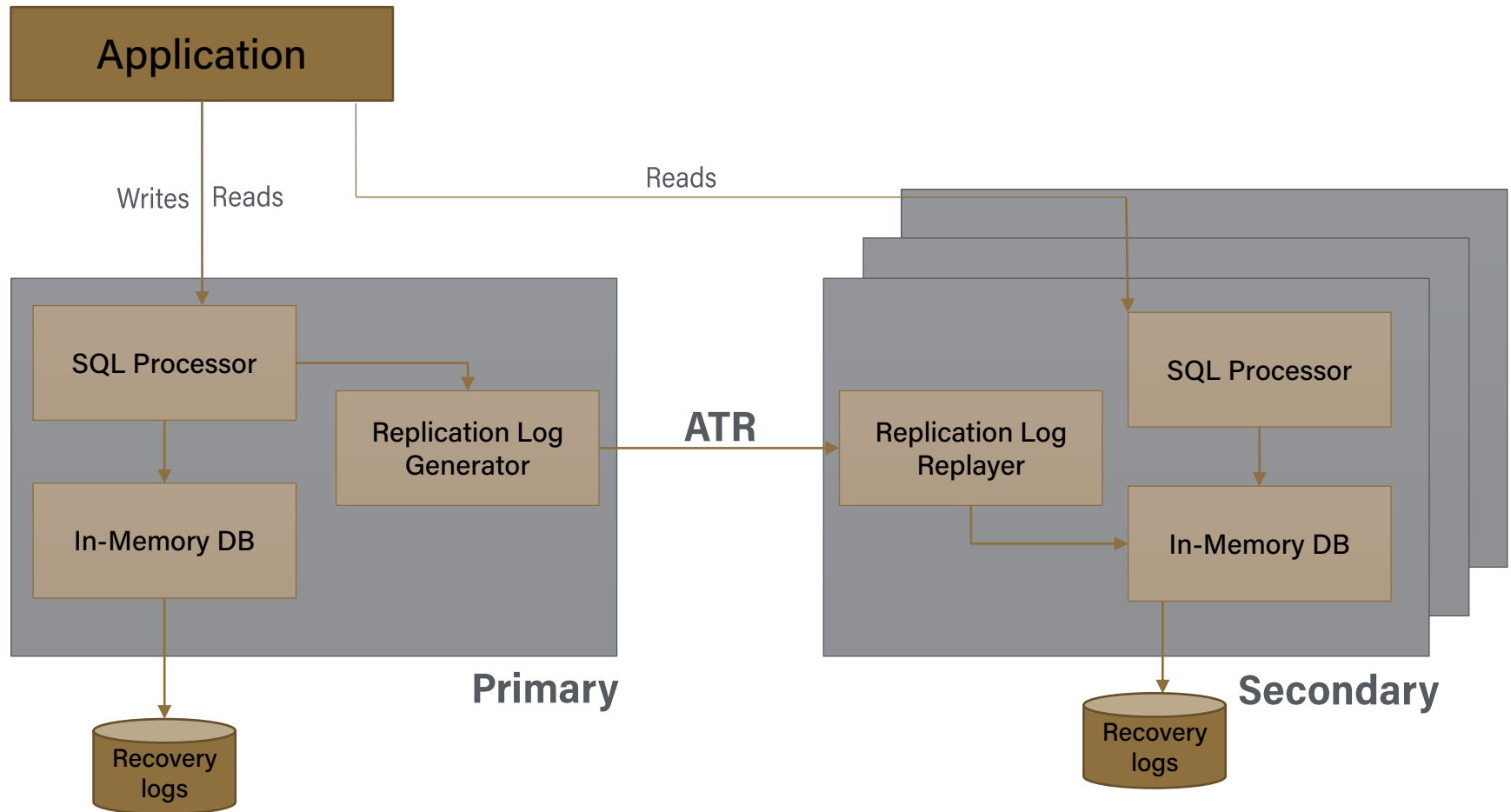
# *Handling Node Failures*

- If an Oracle RAC instance fails, then the IMCUs on the failed instance are unavailable.

- Consequently, a query that needs data stored in the inaccessible IMCUs must read it from somewhere else: the database buffer cache, flash storage, disk, or mirrored IMCUs in other IM column stores.

- Duplication helps to provide fault tolerance because if one node fails, then the mirrored columnar data is accessible from a different node.

- Worst case, queries issued against missing data do not fail. Instead, queries access the data either from the database buffer cache or permanent storage, which may just negatively affect performance.

Source: <u>Oracle Database R21 In-Memory Deploying IM Column Stores in Oracle RAC</u>

# SAP HANA

## Disaggl

Presented by
Shubham Pandey

## Architecture

# Design Decisions

In-Database Replication



**Primary**

Source: http://www.vldb.org/pvldb/vol10/p1598-han.pdf

# Design Decisions

Asynchronous Replication



Primary                                    Secondary

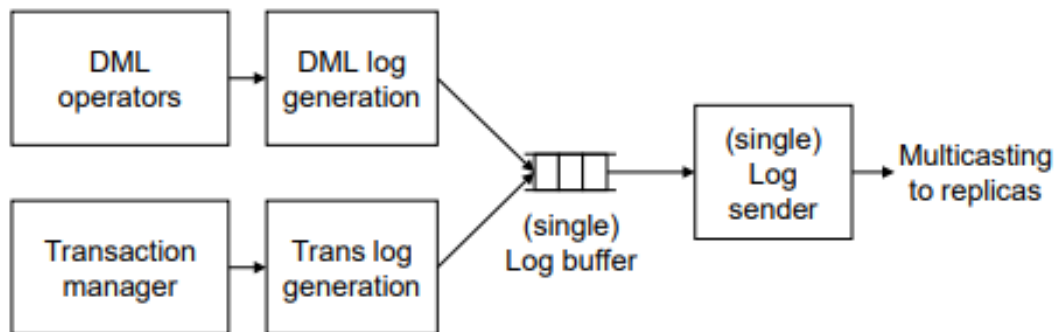Source: http://www.vldb.org/pvldb/vol10/p1598-han.pdf

## Design Decisions

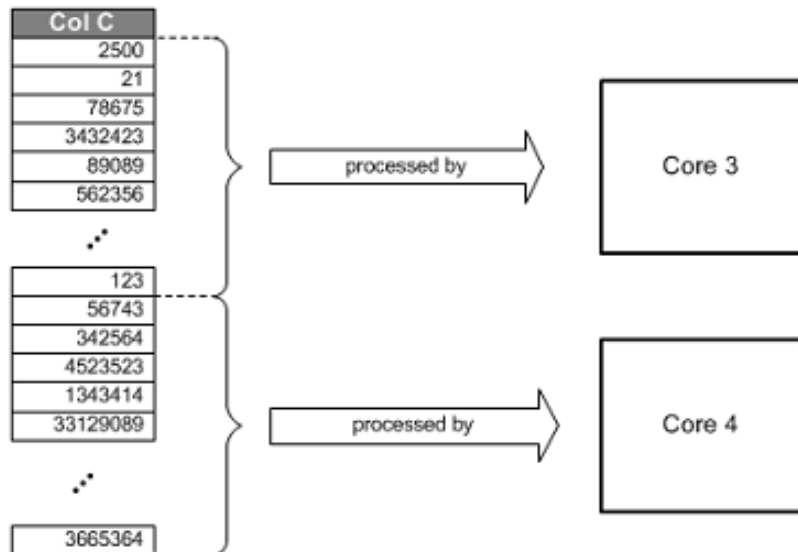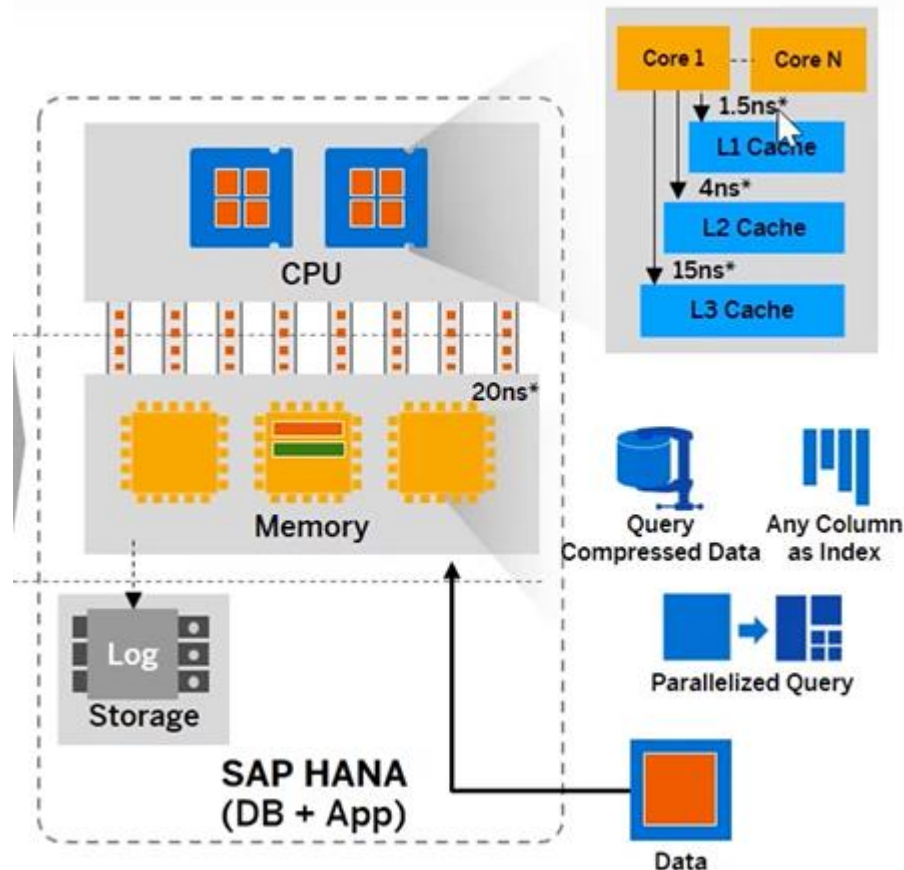Early Log Shipping

# Design Decisions

Log generator and Sender



Source: http://www.vldb.org/pvldb/vol10/p1598-han.pdf

# Design Decisions

Parallel Log Replay



*session_id % n*

# SAP Hana

## Parallel Processing

# Exploiting Modern Hardware

PURDUE UNIVERSITY®

# Fault Tolerance



- Continuous synchronization of HANA DB to secondary location

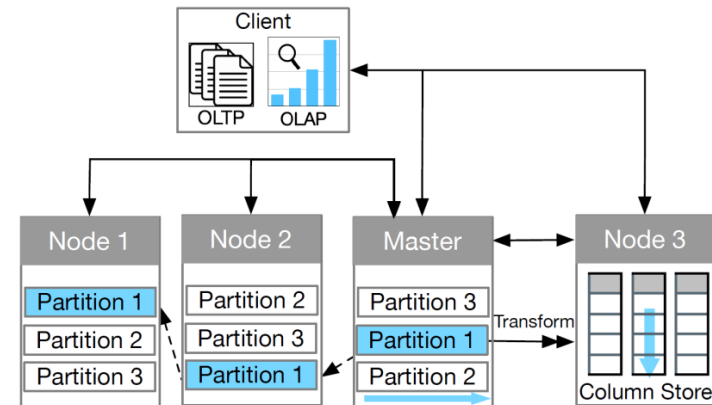- Secondary location can be in same data center or a different one

Source: https://blogs.sap.com/2021/12/31/sap-hana-platform-fault-tolerance-features-high-availability/

**PURDUE UNIVERSITY**

# TiDB: A Raft Based HTAP Database

## Disaggl

Presented by
Satya Sai Bharath Vemula

# System Goals

## System on Higher Level

- Provide Isolation for OLTP and OLAP Queries
  - Providing separate data stores (Row-based and Column-based) for OLTP and OLAP

- Providing Consistent View of Data / Real time data synchronization between the two data stores by extending using state machine-based consensus algorithm (*RAFT*)

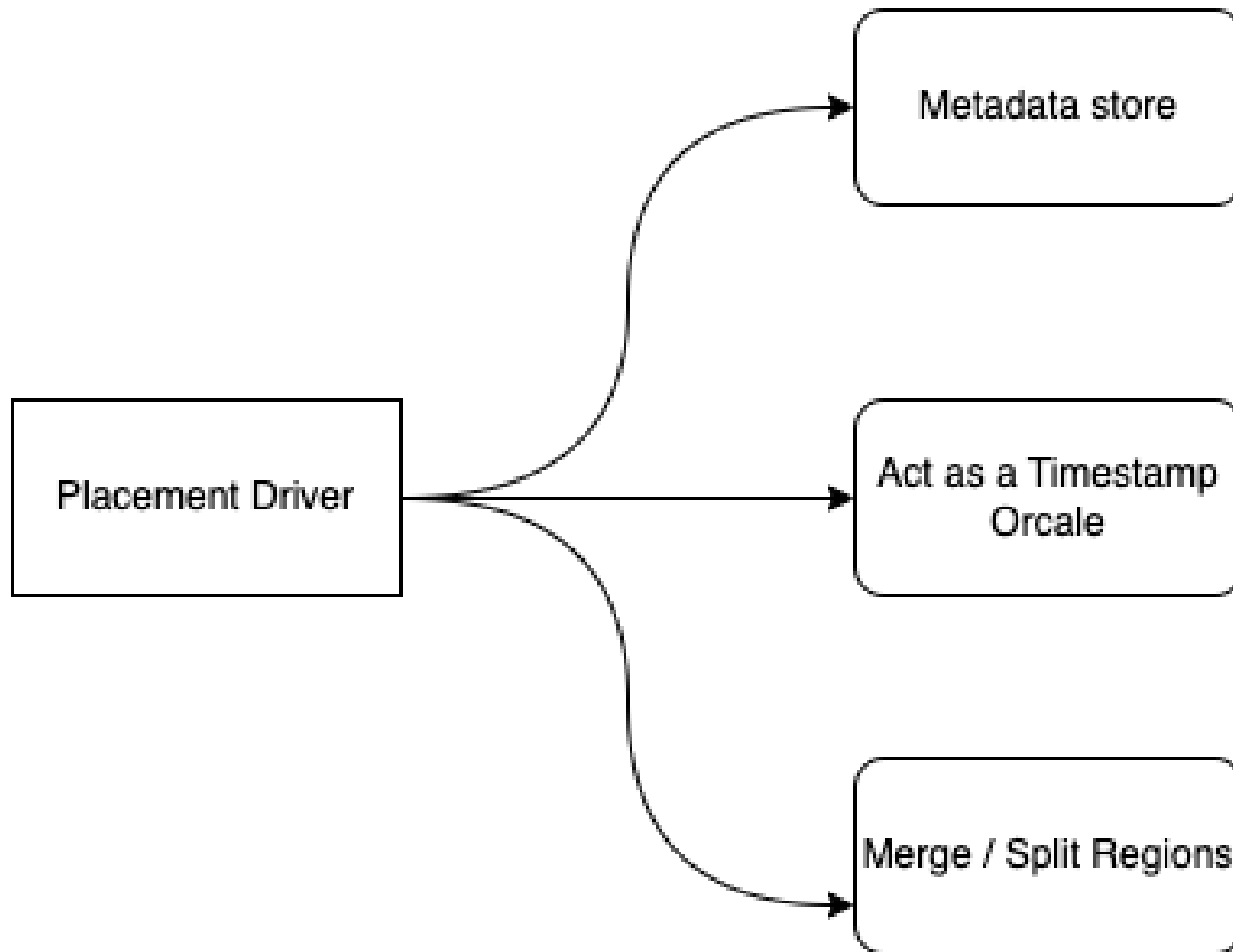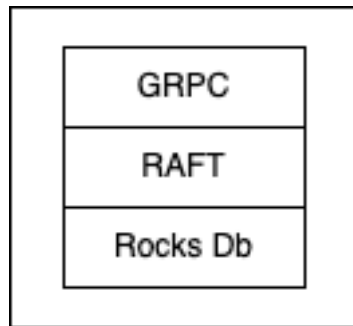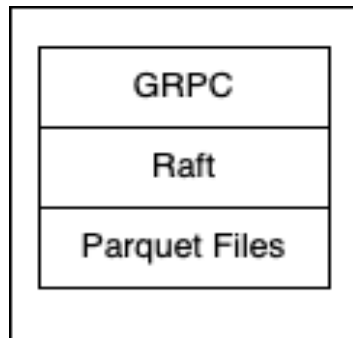- Highly Scalable and Available System with efficient query processing



Source: Guoliang Li and Chao Zhang. 2022. HTAP Databases: What is New and What is Next. (SIGMOD '22)

2/23/2023 | 38

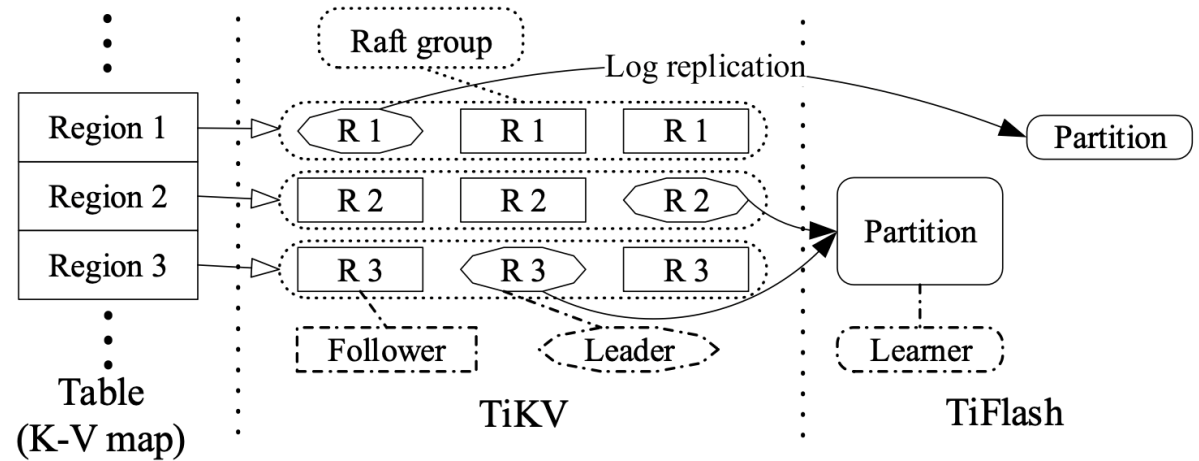# TiDB Architecture

# Availability / Failure Tolerance

Failure Tolerance in both TikV and TiFlash is maintained through replicas

- For each TiKV region if there are 2*m + 1 nodes there will be no observed down time if m + 1 nodes are functional in the same network partition.

  - This Property is because of RAFT Availability Semantics

  - OLTP queries for a region will not be functional if more than m + 1 nodes in a region have failed.

- TiFlash supports multiple replicas for the same data

  - If there are N replicas for a single TiFlash Node, The system can bear a failure of N –1 replicas

  - OLAP queries can still be functional if all the TiFlash Nodes are down as the system can also read data from TiKV nodes for serving analytical queries.
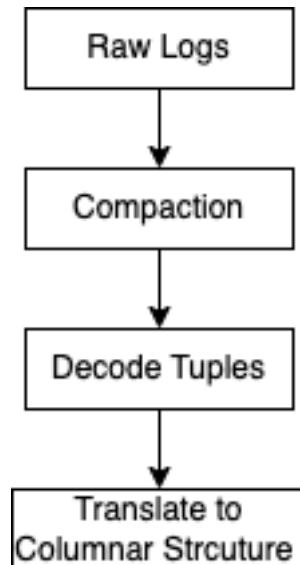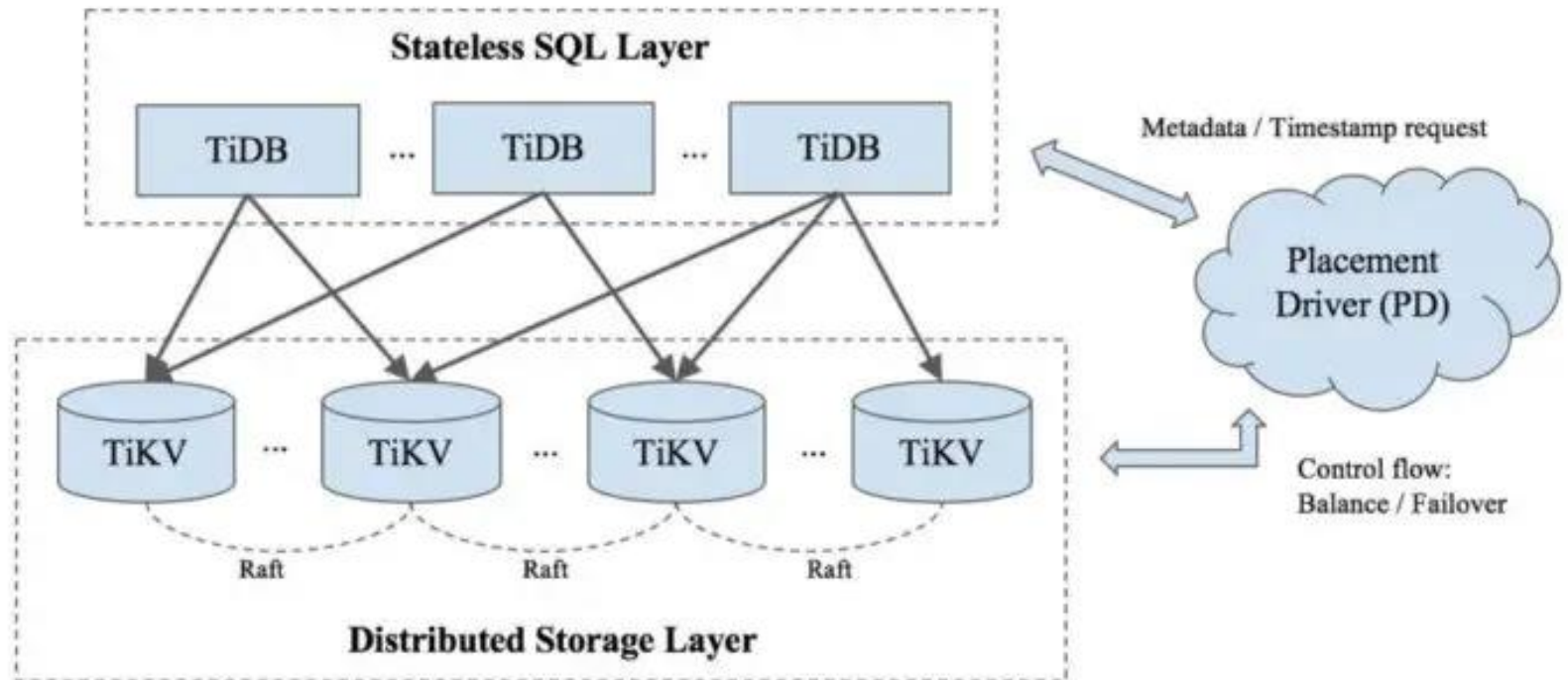
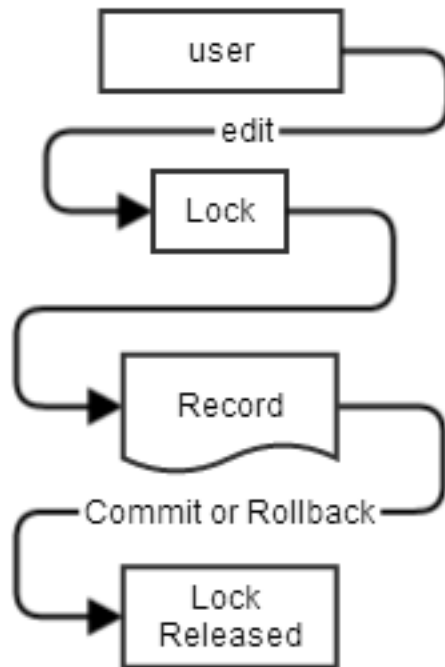**Table 1: Log replaying and decoding**



| | |
|---|---|
| Raw logs | {1}{insert}{prewritten@1}{k1→ (a1, b1)} <br> {2}{insert}{prewritten@2}{k2→ (a2, b2)} <br> {3}{update}{prewritten@3}{k3→ (a3, b3)} <br> {1}{insert}{rollbacked@1} <br> {2}{insert}{committed#4} <br> {3}{update}{committed#5} <br> {4}{delete}{prewritten@6}{k4} <br> {4}{delete}{committed#7} |
| Compacted logs | {2}{insert}{prewritten@2}{k2→ (a2, b2)} <br> {3}{update}{prewritten@3}{k3→ (a3, b3)} <br> {2}{insert}{committed#4} <br> {3}{update}{committed#5} <br> {4}{delete}{prewritten@6}{k4} <br> {4}{delete}{committed#7} |
| Decoded tuples | {insert}{#4}{k2→ (a2, b2)} <br> {update}{#5}{k3→ (a3, b3)} <br> {delete}{#7}{k4} |
| Columnar data | {insert,update,delete,} <br> {#4,#5,#7,} <br> {k2,k3,k4,} <br> {a2,a3,,} <br> {b2,b3,,} |

Source: https://www.vldb.org/pvldb/vol13/p3072-huang.pdf

# OLTP Distributed Transaction



www.adfjavacodes.blogspot.in

Pessimistic Locking

Optimistic Locking

PURDUE
UNIVERSITY®

Source: https://docs.pingcap.com/tidb/dev/tispark-overview

# Project Plan

- **Study HTAP systems with diverse architectures (e.g. data organisation and synchronization)**

  - Compare and Contrast such systems based on their design decisions

  - Assess the strengths and weakness of these systems (in terms of AP/TP throughput, scalability, among others)

- **Study the systems that have evolved from OLTP to HTAP (e.g. MemSQL, IBM dashDB, and more)**

# APPENDIX
# DETAILED DESCRIPTION

PURDUE UNIVERSITY®

# SAP Hana

## Replay Algorithms – DML log entry

**Require:** A DML log entry $L$.
1: Find the transaction object $T$ for $L.TransactionID$.
2: **if** $T$ is empty **then**
3:     Create a transaction object for $L.TransactionID$.
4: **end if**
5: **if** $L.OperationType = Insert$ **then**
6:     Insert $L.Data$ into the table $L.\tau$.
7:     Set the inserted record's $RVID$ as $L.\alpha$.
8: **else if** $L.OperationType = Delete$ **then**
9:     **while** true **do**
10:         Find the record version $R$ whose $RVID$ equals
11:     to $L.\beta$ in the table $L.\tau$.
12:         **if** $R$ is not empty **then**
13:             Delete $R$. **return**
14:         **end if**
15:     **end while**
16: **else if** $L.OperationType = Update$ **then**
17:     **while** true **do**
18:         Find the record version $R$ whose $RVID$ equals
19:     to $L.\beta$ in the table $L.\tau$.
20:         **if** $R$ is not empty **then**
21:             Update $R$ with $L.Data$ and $L.\alpha$. **return**
22:         **end if**
23:     **end while**
24: **end if**

$\tau - Table\ ID$

$\beta - Before\ update\ RVID$

$\alpha - $ After update RVID

Source: http://www.vldb.org/pvldb/vol10/p1598-han.pdf

## Replay Algorithms

**Replay a precommit log entry**

**Require:** A precommit log entry $L$.
1: Find the transaction object $T$ for $L.TransactionID$.
2: Mark $T$'s state as *precommitted*.

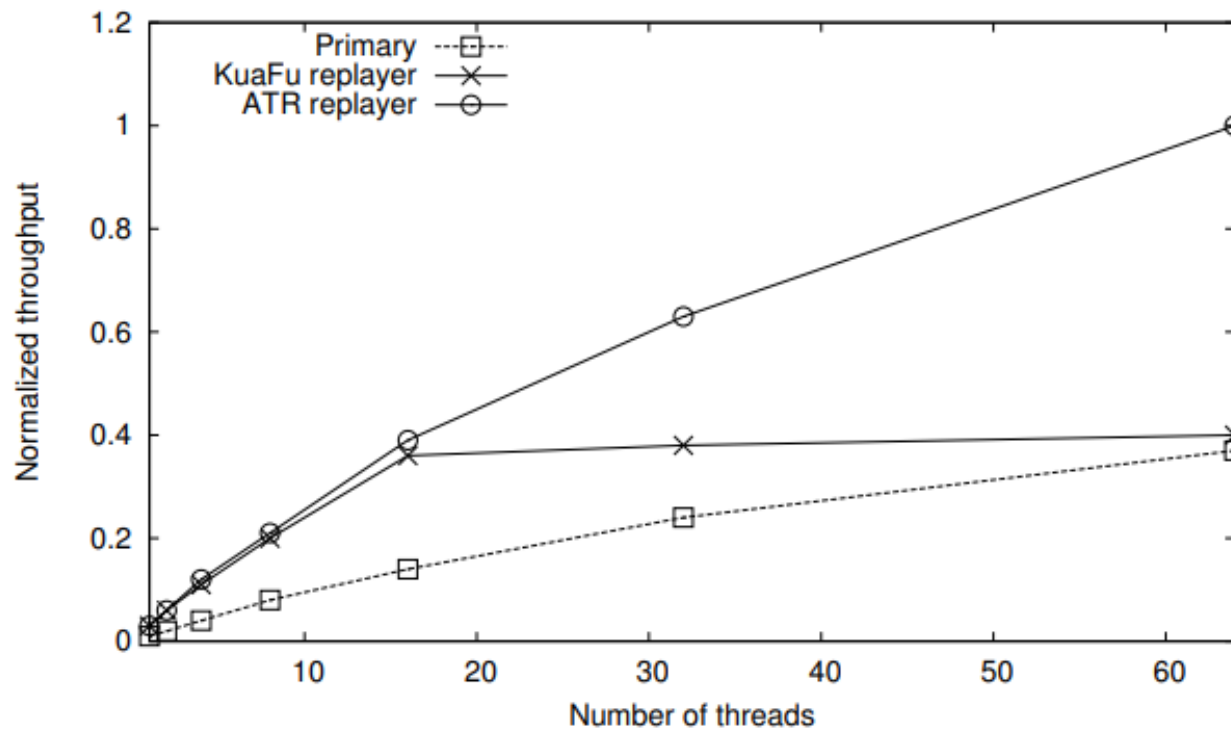**Replay an abort log entry**

**Require:** An abort log entry $L$.
1: Find the transaction object $T$ for $L.TransactionID$.
2: Abort $T$.

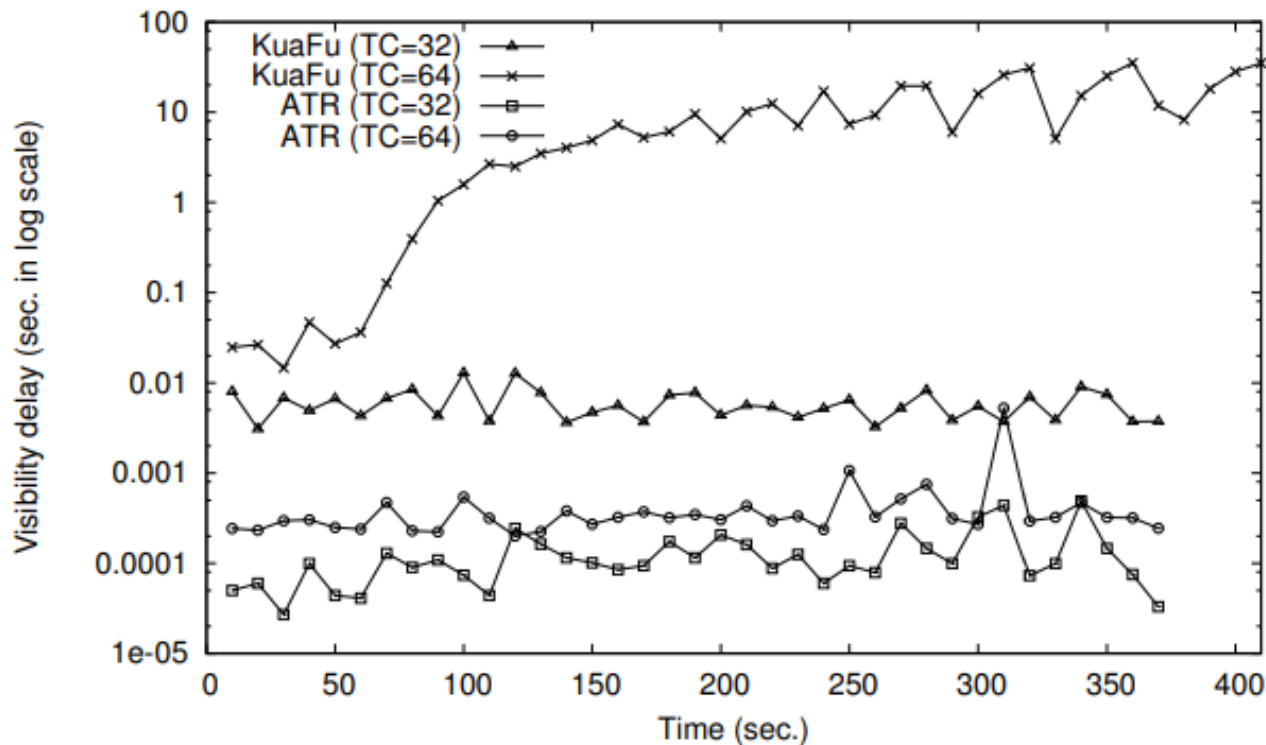**Replay a commit log entry**

**Require:** A commit log entry $L$.
1: Find the transaction object $T$ for $L.TransactionID$.
2: Wait until $T$'s state becomes *precommitted*.
3: Increment the transaction commit timestamp of the replica server by marking the $T$'s generated record versions with a new commit timestamp value.

Source: http://www.vldb.org/pvldb/vol10/p1598-han.pdf

**PURDUE UNIVERSITY**

# Multi-core Scalability with Parallel Log Replay



**TPC-C benchmark**

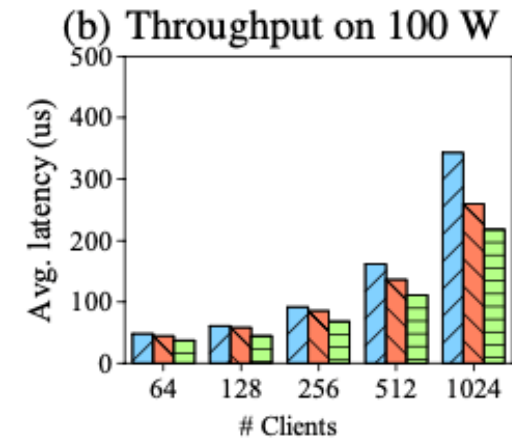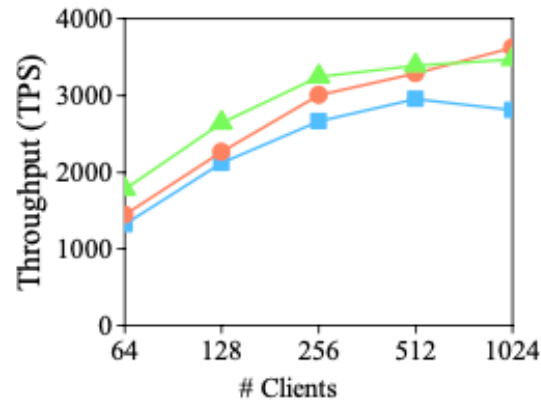Source: http://www.vldb.org/pvldb/vol10/p1598-han.pdf

## Visibility Delay



**Visibility delay at replica while running TPC-C benchmark at primary**

Source: http://www.vldb.org/pvldb/vol10/p1598-han.pdf

# Multi Replica scalability under mixed (OLTP+OLAP) workload



Source: http://www.vldb.org/pvldb/vol10/p1598-han.pdf

(a) Throughput on 50 W

(b) Throughput on 100 W

(c) Throughput on 200 W

(d) Latency on 200 W

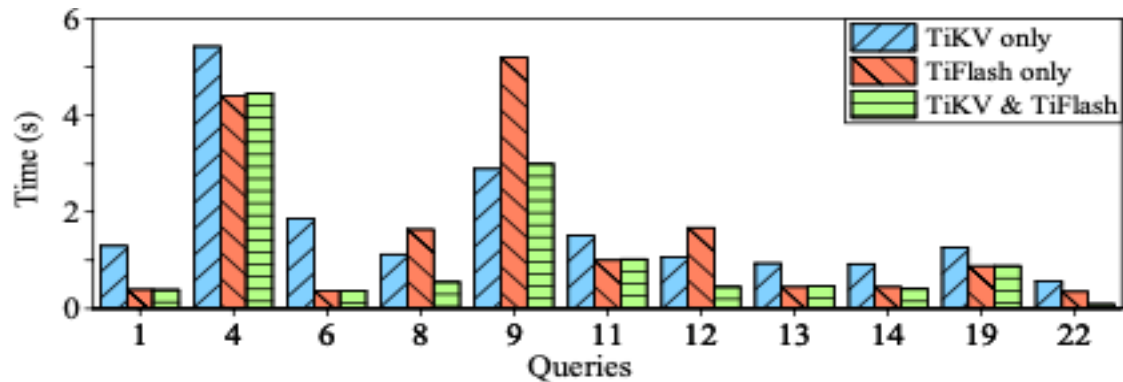Source: https://www.geeksforgeeks.org/three-phase-commit-protocol/

**Figure 8: Choice of TiKV or TiFlash for analytical queries**



**Figure 9: Performance comparison of CH-benCHmark analytical queries**

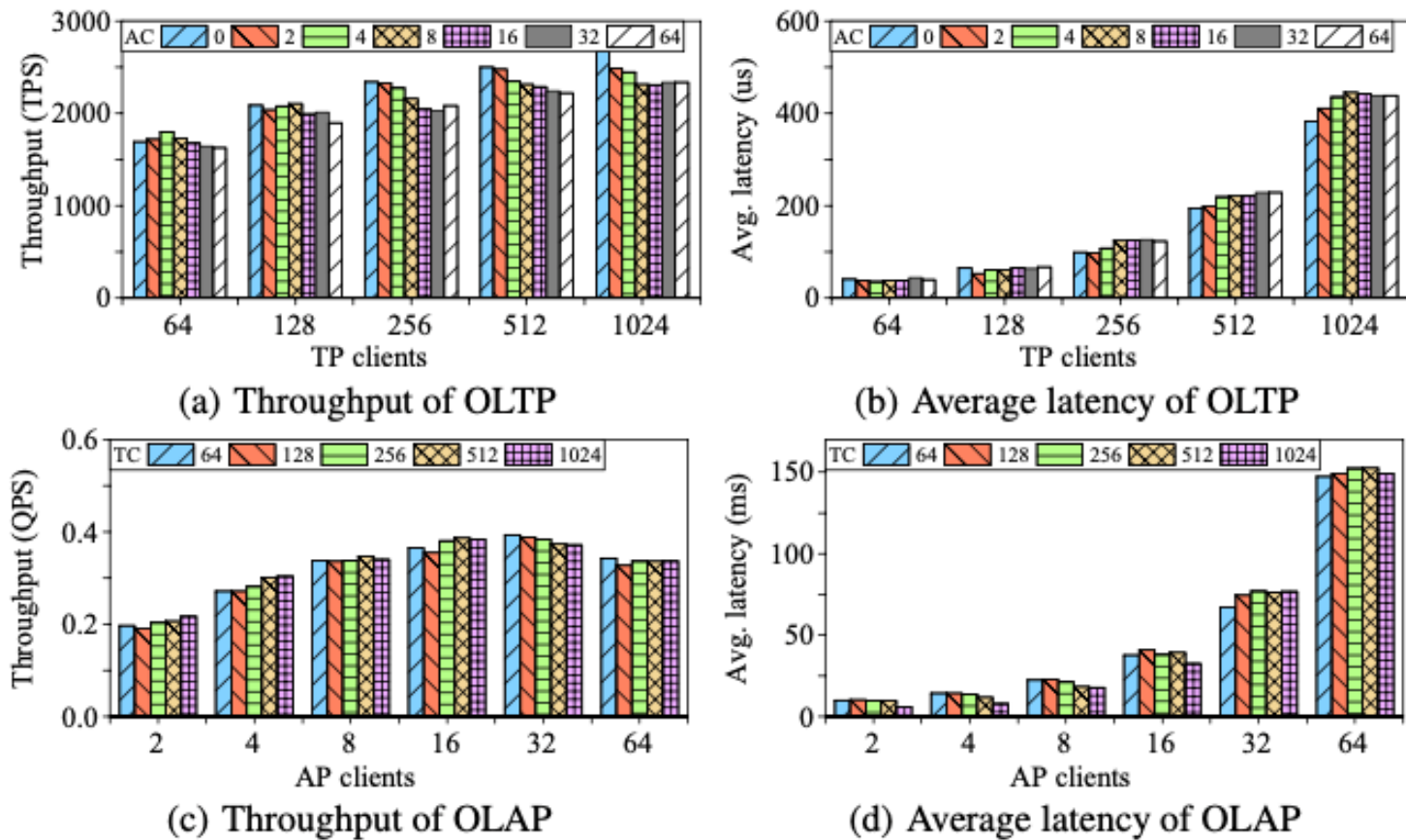Source: https://www.vldb.org/pvldb/vol13/p3072-huang.pdf

Figure 10: HTAP performance of TiDB

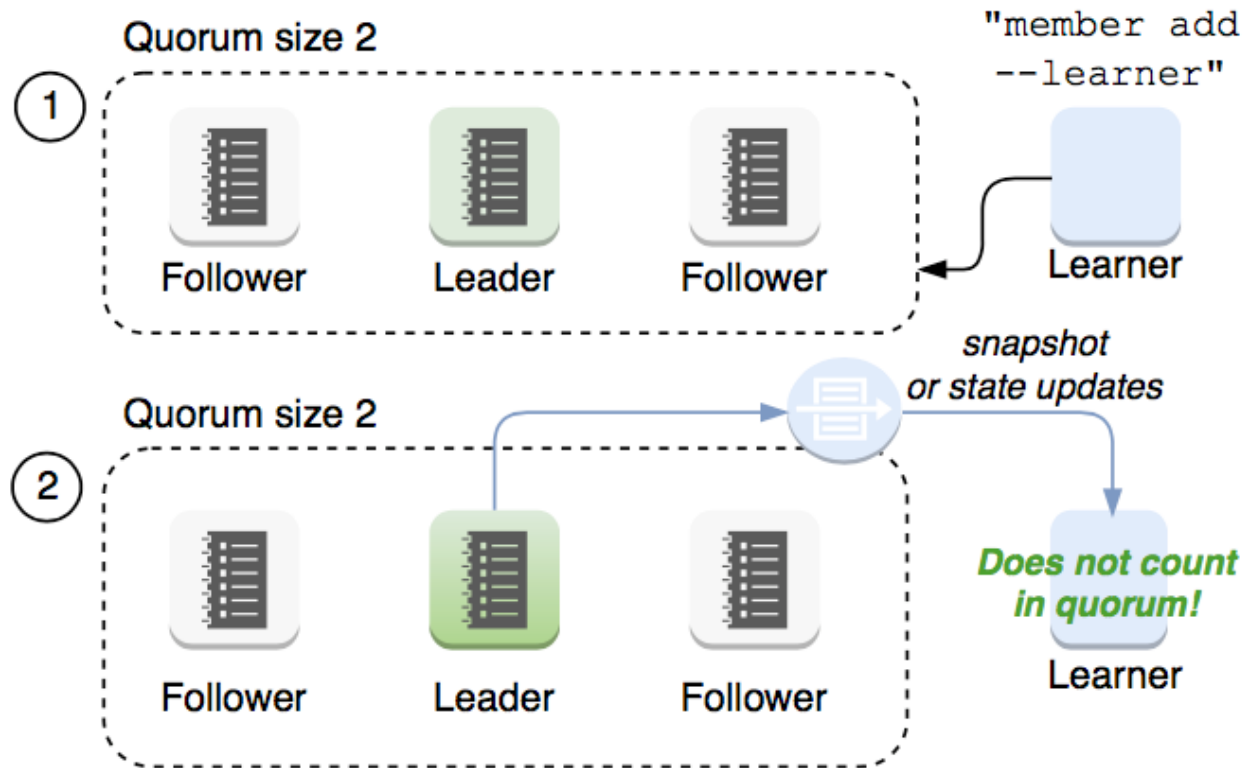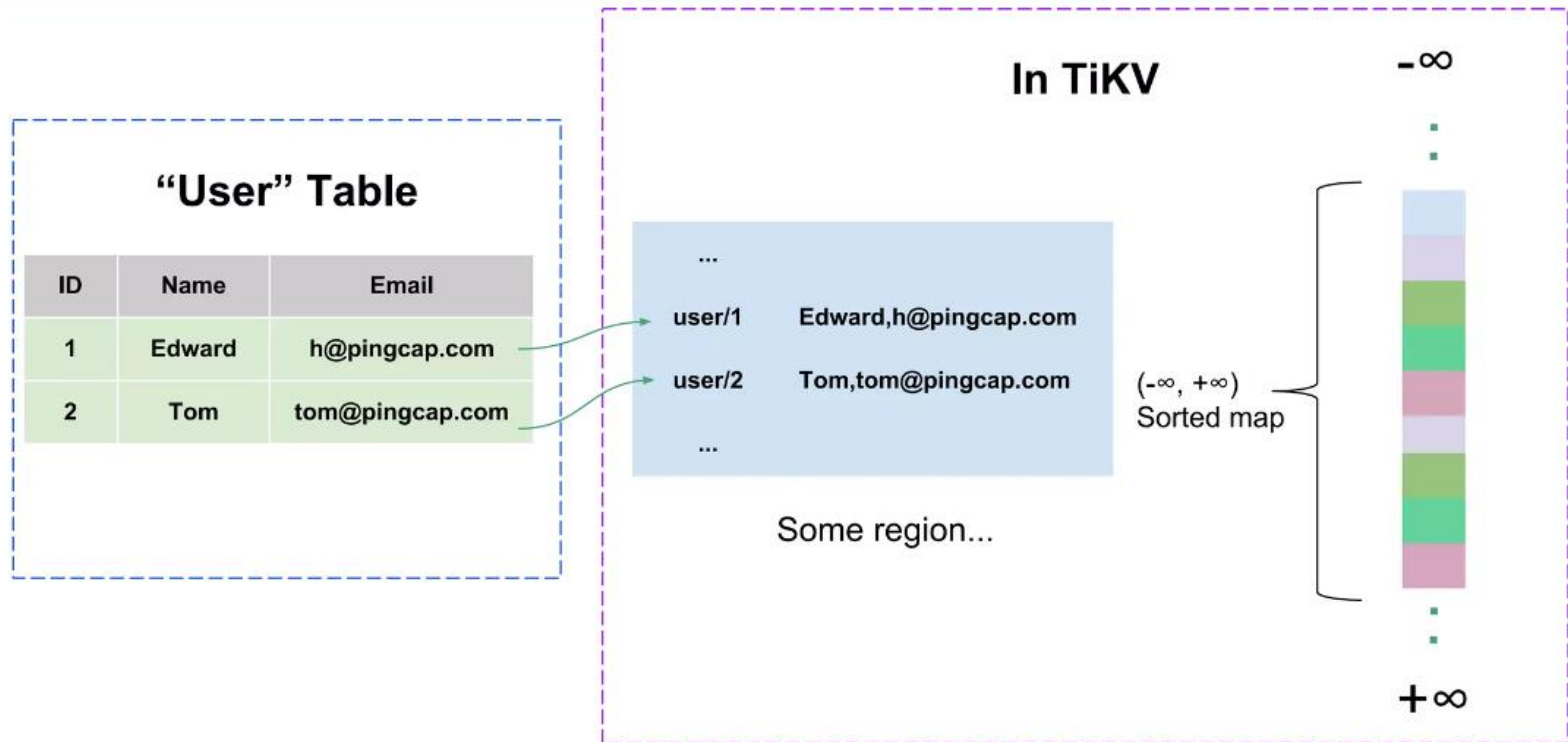Source: https://www.vldb.org/pvldb/vol13/p3072-huang.pdf

Figure 10. Add a learner node as a non-voting member. Wait until learner node catches up to leader's logs. Until then, learner node neither votes nor counts towards quorum.

Figure 4: The columnar delta tree

$$C_{\text{opt\_scan}} = \min(C_{\text{col\_scan}}, C_{\text{row\_scan}}, C_{\text{index\_scan}})$$

$$C_{\text{row\_scan}} = S_{\text{tuple}} \cdot N_{\text{tuple}} \cdot f_{\text{scan}} + N_{\text{reg}} \cdot f_{\text{seek}}$$

$$C_{\text{col\_scan}} = \sum_{j=1}^{m} \left( S_{\text{col\_j}} \cdot N_{\text{tuple}} \cdot f_{\text{scan}} + N_{\text{reg\_j}} \cdot f_{\text{seek}} \right)$$

$$C_{\text{index\_scan}} = S_{\text{index}} \cdot N_{\text{tuple}} \cdot f_{\text{scan}} + N_{\text{reg}} \cdot f_{\text{seek}} + C_{\text{double\_read}}$$

$$C_{\text{double\_read}} = \begin{cases} 0 & \text{(if without double read)} \\ S_{\text{tuple}} \cdot N_{\text{tuple}} \cdot f_{\text{scan}} + N_{\text{tuple}} \cdot f_{\text{seek}} \end{cases}$$

C -> Cost
S -> Size
F -> Seek / Scan Cost
N -> Number
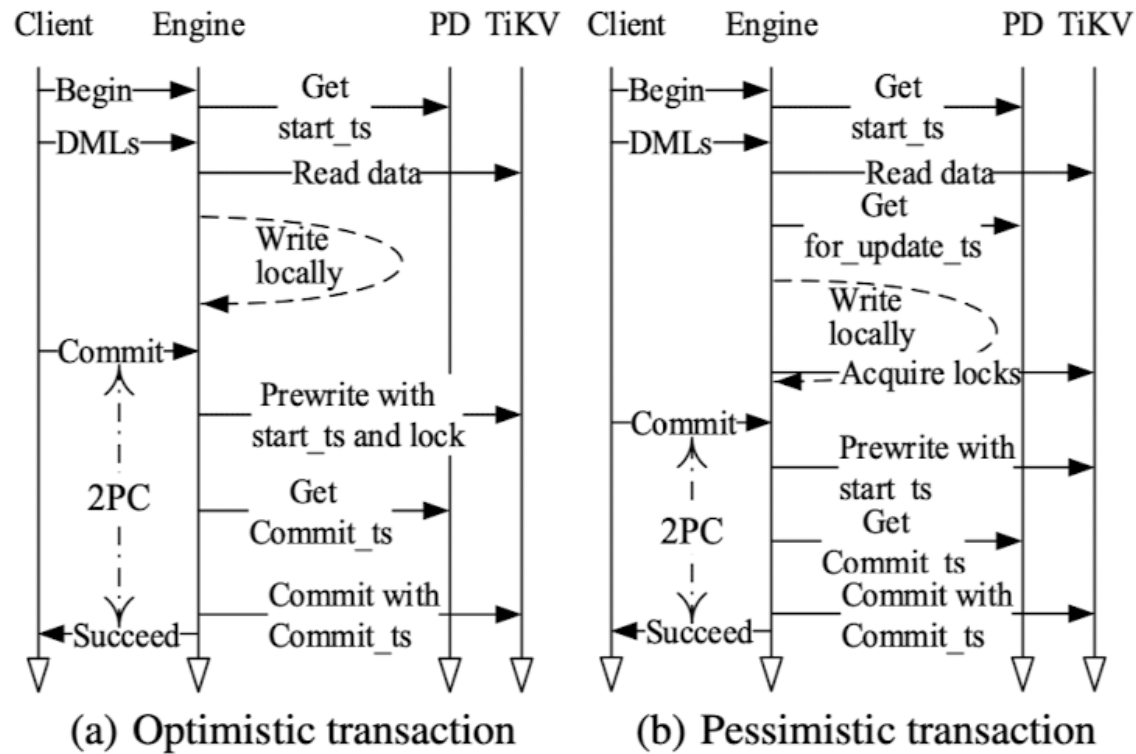
Source: https://www.vldb.org/pvldb/vol13/p3072-huang.pdf

**Figure 5: The process of optimistic and pessimistic transaction**