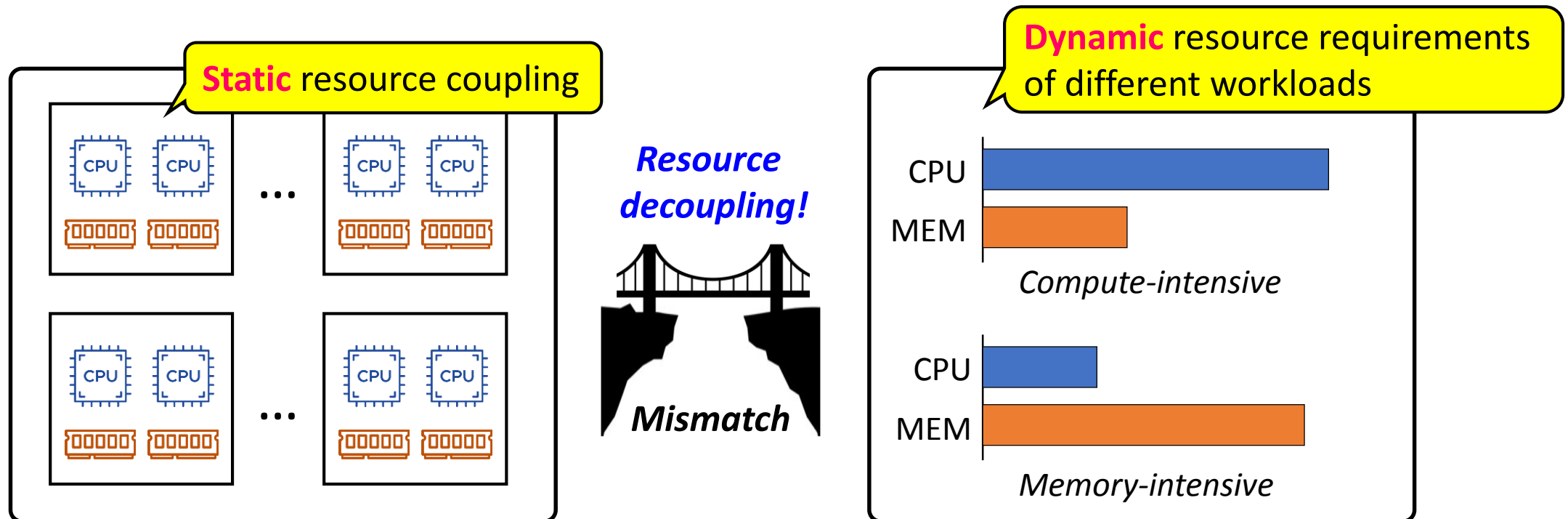# Motor: Enabling Multi-Versioning for Distributed Transactions on Disaggregated Memory

**Ming Zhang**, Yu Hua, Zhijun Yang

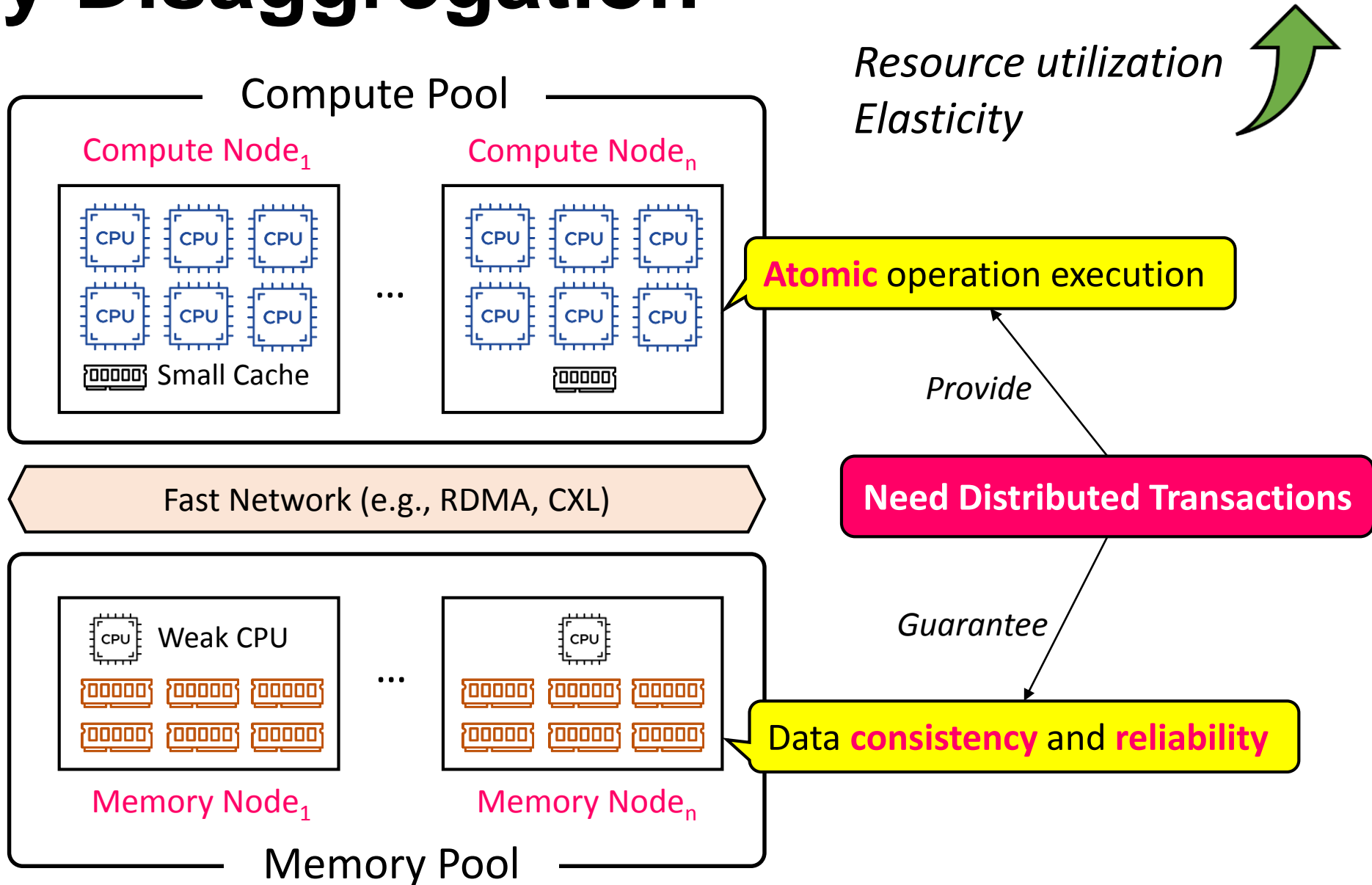*Huazhong University of Science and Technology, China*

# Low Memory Utilization in the Cloud

➢ Below ~60% [1-4]

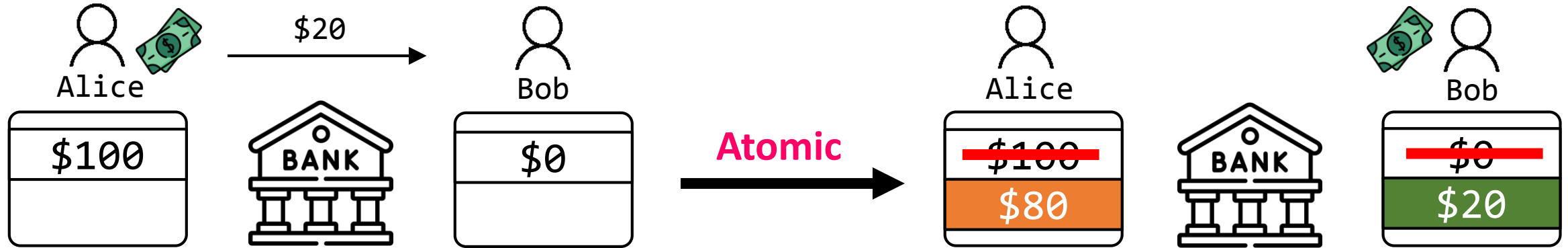➢ One major reason: monolithic server

**Static** resource coupling

*Resource decoupling!*

*Mismatch*

**Dynamic** resource requirements of different workloads

CPU
MEM
*Compute-intensive*

CPU
MEM
*Memory-intensive*

[1] MemTrade@SIGMETRICS'23, Borg@EuroSys'20, LegoOS@OSDI'18
[2] Google Production Cluster Trace. https://github.com/google/cluster-data
[3] Alibaba Production Cluster Trace. https://github.com/alibaba/clusterdata
[4] Snowflake Dataset. https://github.com/resource-disaggregation/snowset

# Memory Disaggregation

# Transcription



```
Txn begin
    Alice: $100 -> $80
      Bob:   $0 -> $20
Txn end
```

*Transaction*

# Distributed Transaction



- Concurrency control
- Atomic commit
- Replication

Alice

$100

BANK

Bob

$0

**Atomic**

Alice

$100
$80

BANK

Bob

$0
$20

**Coordinator**

Alice
**(Primary)**

**(Backups)**

Bob
**(Primary)**

**(Backups)**

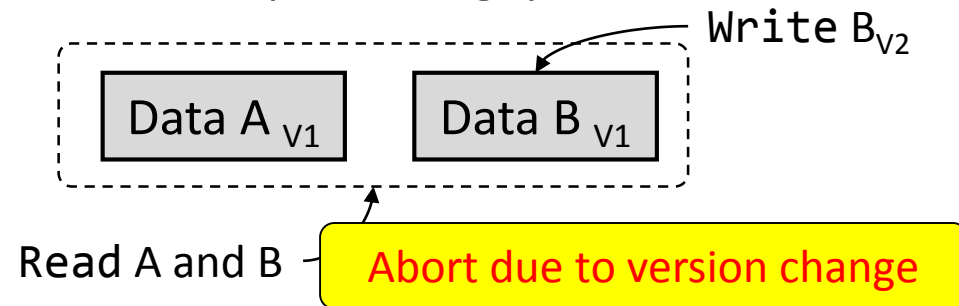*Network*

5

# Execution Model on Disaggregated Memory

# State-of-The-Art Studies

> Single-versioning txn system
  - Based on disaggregated memory[1]

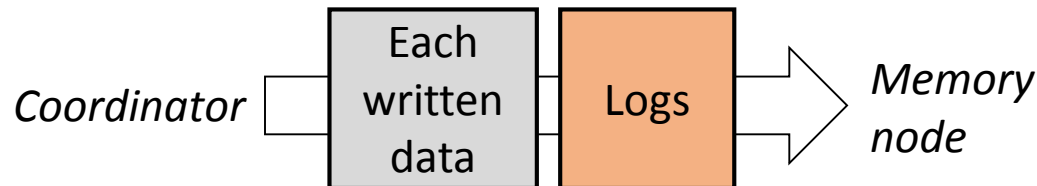☹ Write interrupts read
  - Hamper throughput

Write $B_{V2}$

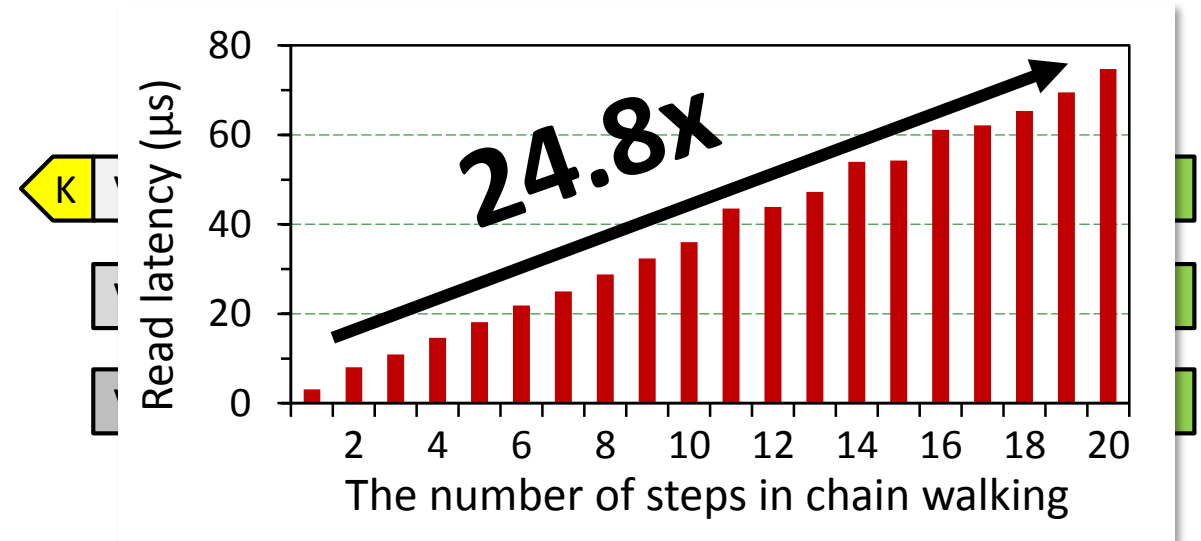| Data A $_{V1}$ | Data B $_{V1}$ |

Read A and B

> Abort due to version change

☹ Substantial write-ahead logs
  - Consume network resources

Coordinator — | Each written data | Logs | ⟹ *Memory node*

> Multi-versioning txn systems
  - Based on monolithic architecture

☹ Inefficient linked version chain



24.8x

Read latency (μs)

The number of steps in chain walking
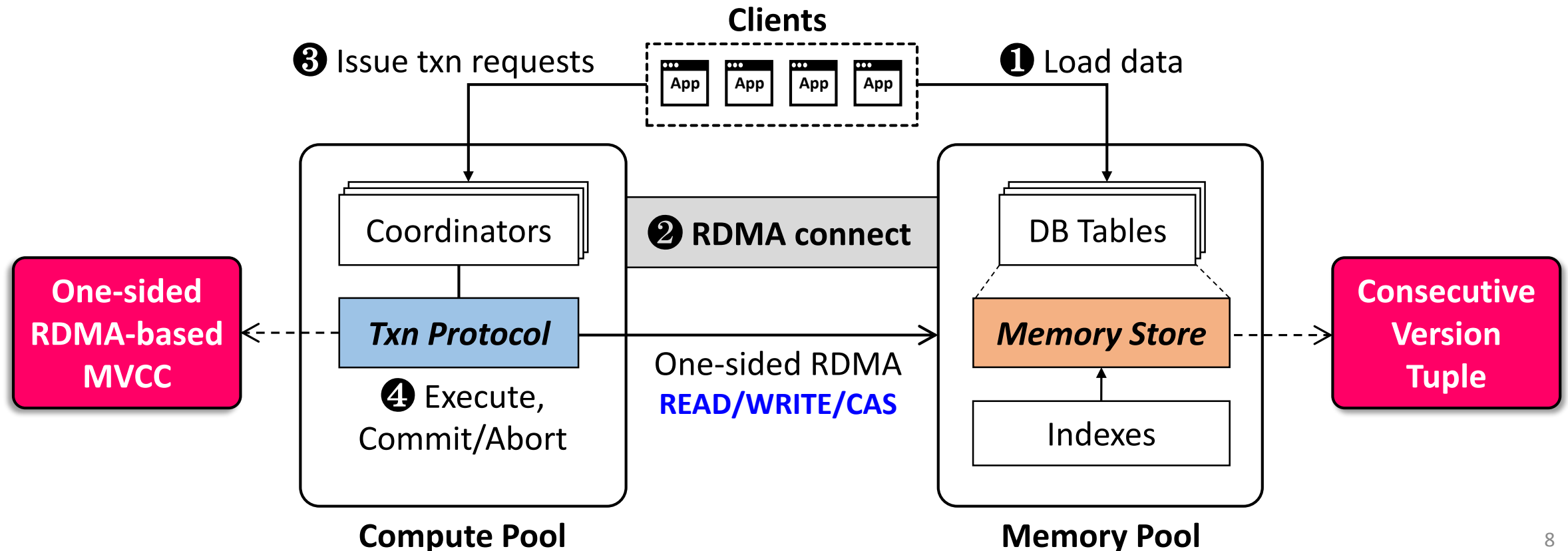
☹ Incompatible transaction protocol
  - Frequent CPU involvement on each data node: timestamp calculation[2], locking[6], validation[7]

[1] FORD@FAST'22    [2] DST@NSDI'21    [3] Postgres    [4] Hekaton@SIGMOD'13    [5] Aurogon@FAST'22
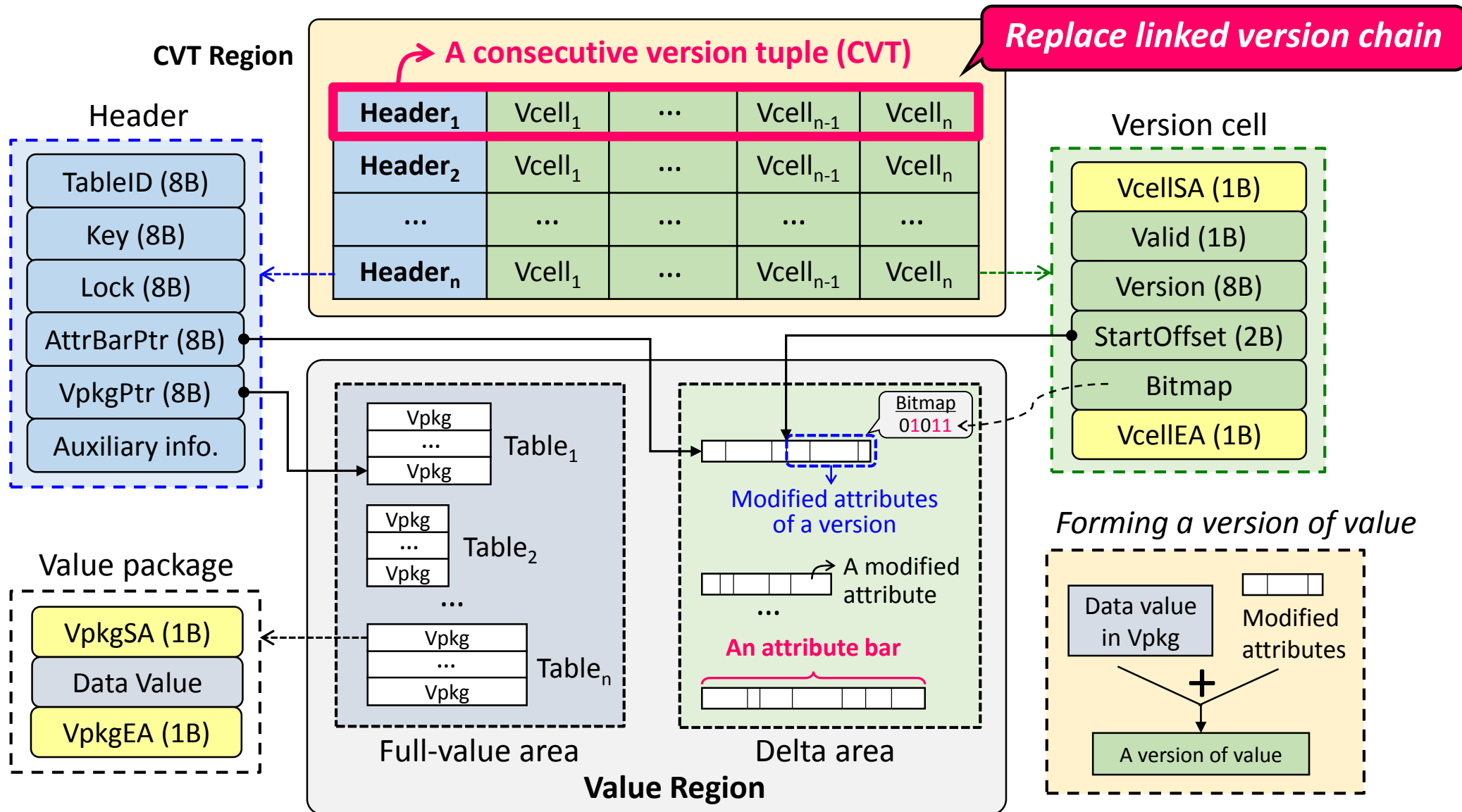[6] FaRMv2@SIGMOD'19    [7] Neumann et al.@SIGMOD'15    [8] MySQL    [9] NAM-DB@VLDB'17
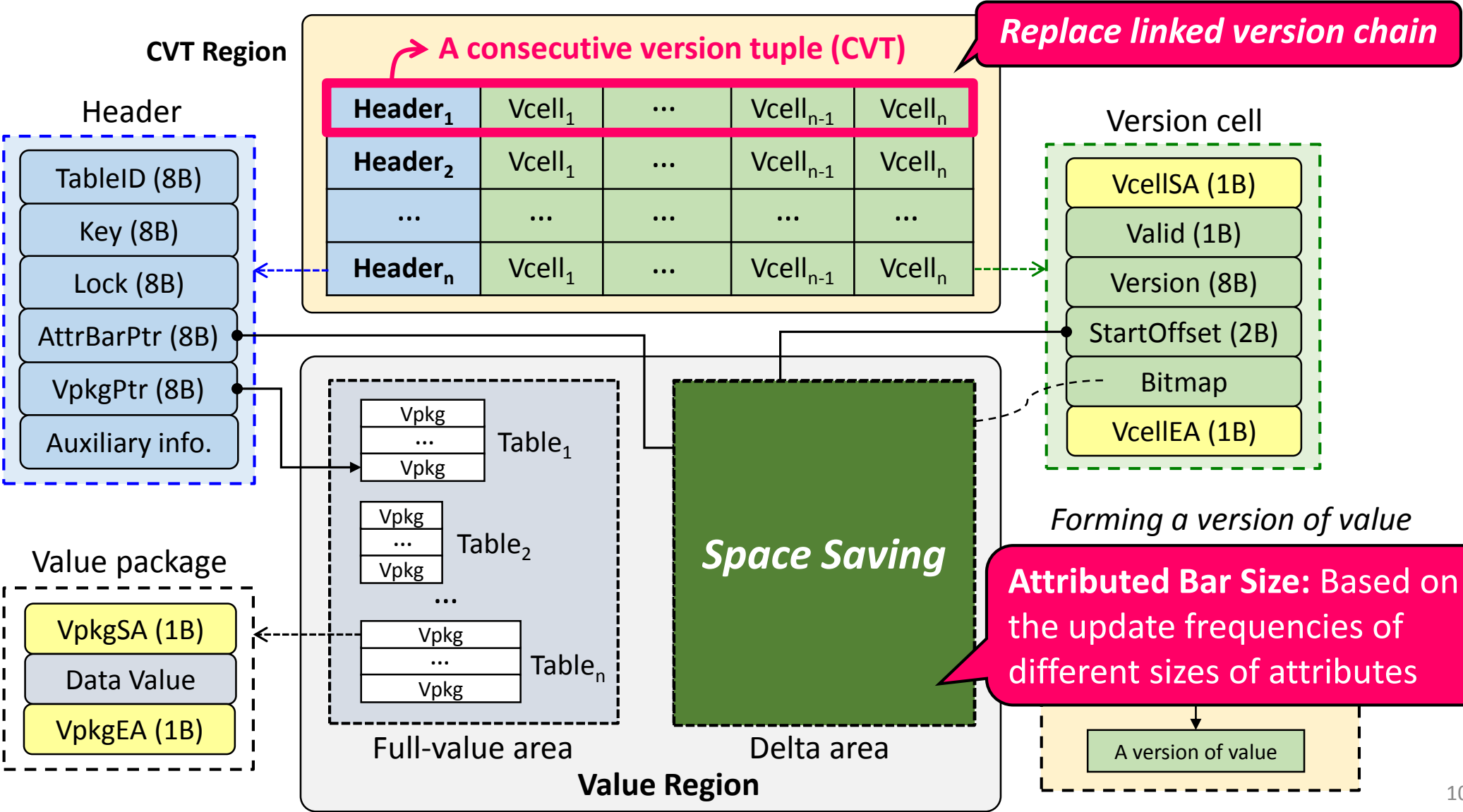
# Motor: Overview

➢ Holistic new design working in harmony
- *Version structure* in memory pool
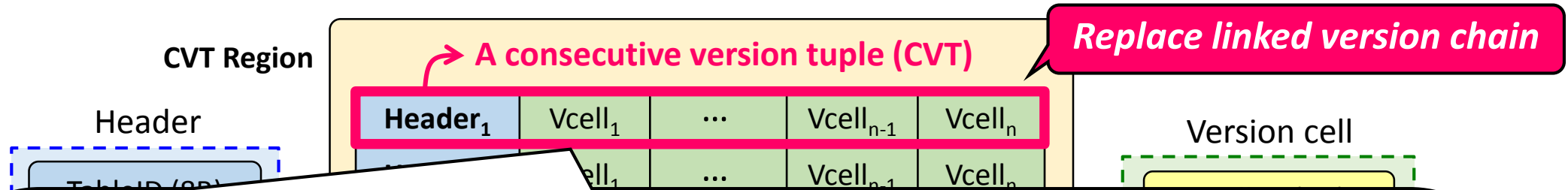- *MVCC protocol* in compute pool
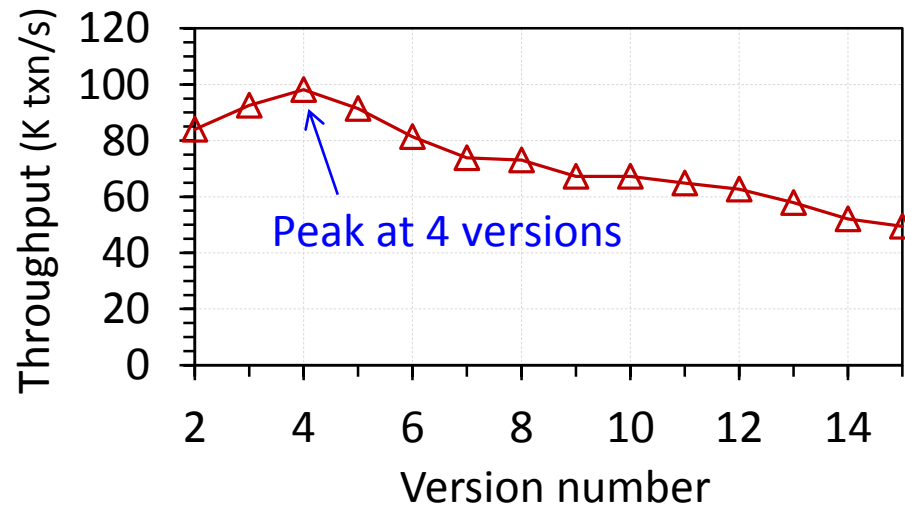
# Consecutive Version Tuple

# Consecutive Version Tuple

CVT Region

**A consecutive version tuple (CVT)**

*Replace linked version chain*

| Header$_1$ | Vcell$_1$ | ... | Vcell$_{n-1}$ | Vcell$_n$ |
|---|---|---|---|---|
| Header$_2$ | Vcell$_1$ | ... | Vcell$_{n-1}$ | Vcell$_n$ |
| ... | ... | ... | ... | ... |
| Header$_n$ | Vcell$_1$ | ... | Vcell$_{n-1}$ | Vcell$_n$ |

Header

- TableID (8B)
- Key (8B)
- Lock (8B)
- AttrBarPtr (8B)
- VpkgPtr (8B)
- Auxiliary info.

Version cell

- VcellSA (1B)
- Valid (1B)
- Version (8B)
- StartOffset (2B)
- Bitmap
- VcellEA (1B)

*Forming a version of value*

Value package

- VpkgSA (1B)
- Data Value
- VpkgEA (1B)

**Value Region**

Full-value area

- Vpkg / ... / Vpkg — Table$_1$
- Vpkg / ... / Vpkg — Table$_2$
- ...
- Vpkg / ... / Vpkg — Table$_n$

*Space Saving*

Delta area

**Attributed Bar Size:** Based on the update frequencies of different sizes of attributes

A version of value

# Consecutive Version Tuple

CVT Region

**A consecutive version tuple (CVT)**

*Replace linked version chain*

Header

TableID (8B)

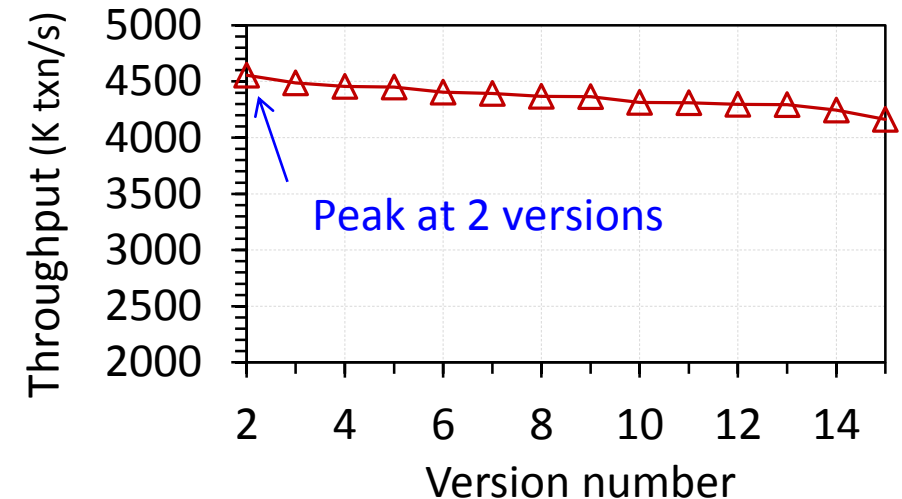| Header$_1$ | Vcell$_1$ | ... | Vcell$_{n-1}$ | Vcell$_n$ |
| ... | ... | Vcell$_{n-1}$ | Vcell$_n$ |

Version cell

**Number of Versions:** depending on `workload characteristics`
- Read-write contention
- Number of accessed records



Peak at 4 versions

TPCC (high contention, long txns)



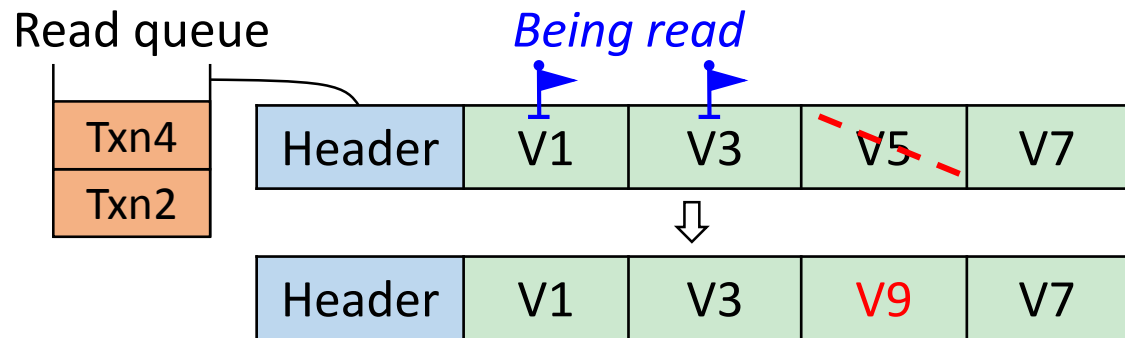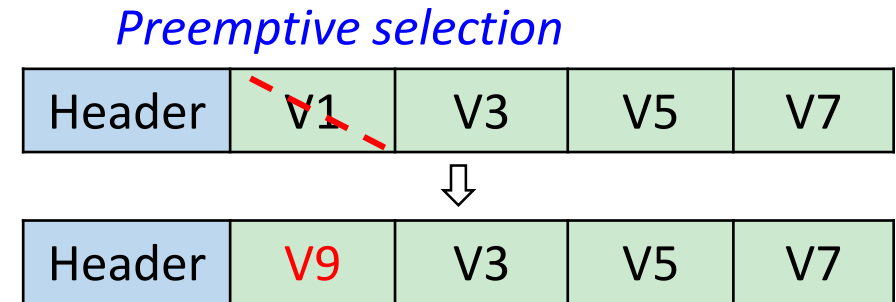Peak at 2 versions

TATP (low contention, short txns)

# Coordinator-Active Garbage Collection

➤ A CVT runs out of space – GC required

➤ Prior systems track transaction states[1-2]

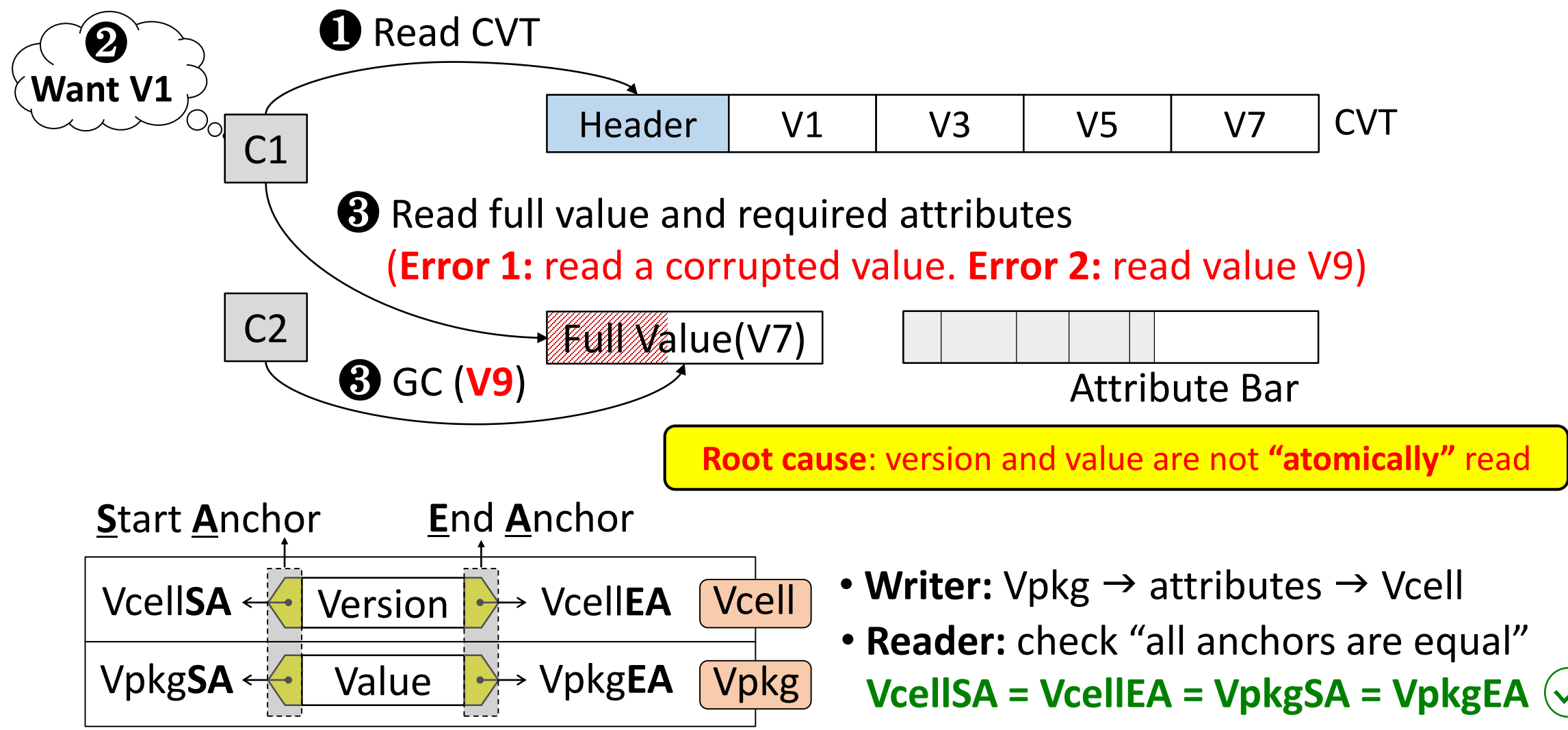- CPU in memory nodes is too weak to frequently track

Read queue      *Being read*        *Preemptive selection*

| Txn4 |
| Txn2 |

| Header | V1 | V3 | V5 | V7 |

⇩

| Header | V1 | V3 | V9 | V7 |

Skip the versions being read

*High overhead for compute nodes to maintain states*

| Header | V1 | V3 | V5 | V7 |

⇩

| Header | V9 | V3 | V5 | V7 |

Overwrite the oldest version

*Simple, no tracking*
*Low abort rate with fast RDMA*

[1] Steam@VLDB'19    [2] FaRMv2@SIGMOD'19

# Anchor-Assisted Read

❷ **Want V1**

❶ Read CVT

C1

| Header | V1 | V3 | V5 | V7 | CVT |
|--------|----|----|----|----|----|

❸ Read full value and required attributes
(**Error 1:** read a corrupted value. **Error 2:** read value V9)

C2

Full Value(V7)

Attribute Bar

❸ GC (**V9**)

**Root cause**: version and value are not **"atomically"** read

**Start Anchor**     **End Anchor**

| VcellSA ← | Version | → VcellEA | Vcell |
|-----------|---------|-----------|-------|
| VpkgSA ← | Value | → VpkgEA | Vpkg |

- **Writer:** Vpkg → attributes → Vcell
- **Reader:** check "all anchors are equal"

**VcellSA = VcellEA = VpkgSA = VpkgEA** ✓

# One-Sided RDMA-Based MVCC



| System | RTTs | Data Version |
|---|---|---|
| DrTM+H@OSDI'18 | 4 | Single Version |
| FORD@FAST'22 | 3 | Single Version |
| FaRMv2@SIGMOD'19 | ≥ 5 | Multiple Version |
| **Motor** | **≤ 4** | **Multiple Version** |

**Deliver MV benefits without long RTTs**

\* Tstart/Tcommit stands for start/commit timestamp

14

# Evaluation

- ## Benchmarks
  - KV store
    - 8B key + 40B value
    - Skewed (skewness tunable)
  - TATP
    - RO/RW: 80%/20%, max 48B
  - SmallBank
    - RO/RW: 15%/85%, 16B
  - TPCC
    - RO/RW: 8%/92%, max 672B
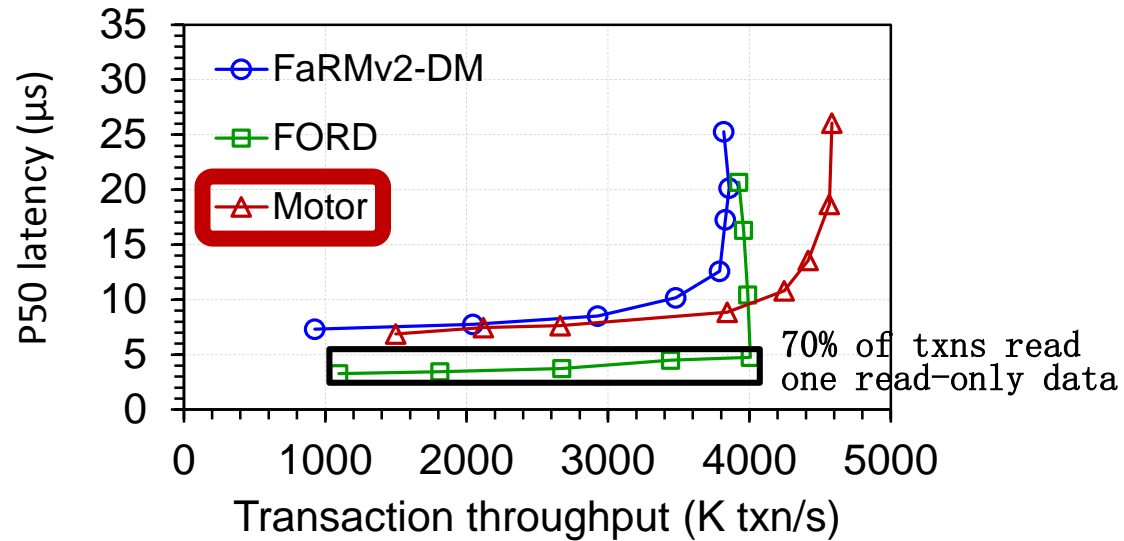- ## Comparisons
  - FaRMv2@SIGMOD'19
  - FORD@FAST'22

Compute Pool

CX-5 RNIC

100 Gbps
InfiniBand switch

CX-5 RNIC          CX-5 RNIC          CX-5 RNIC

Memory Pool

# Performance of Version Structures



Version number = 3, skewness = 0.7



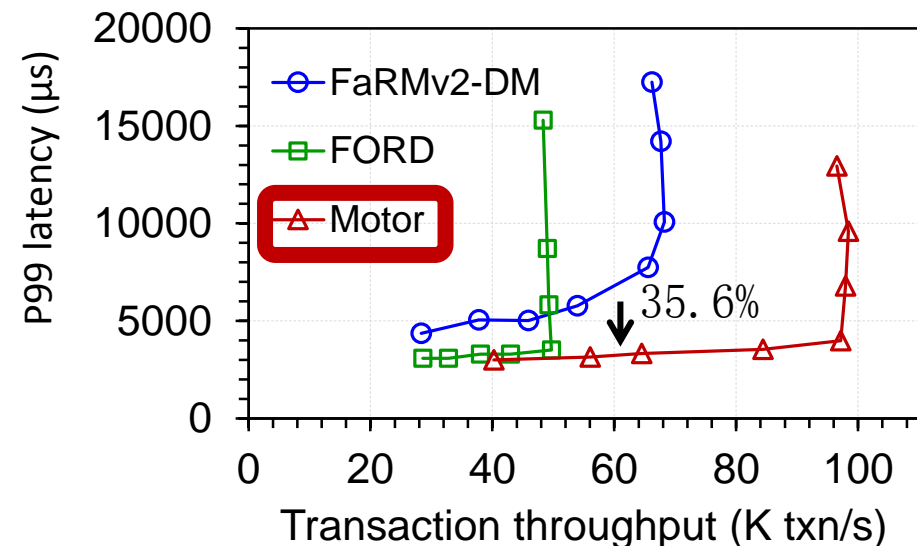Version number = 4, skewness = 0.99
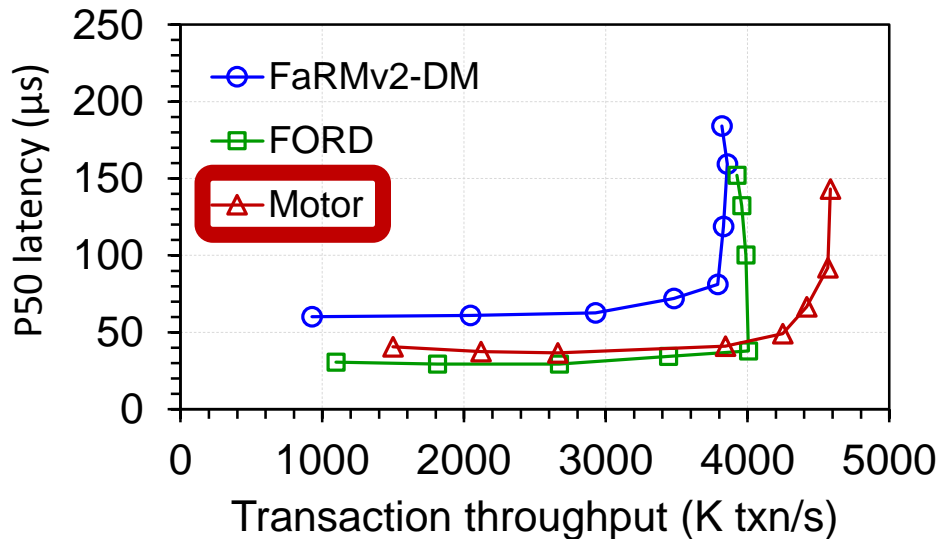
CVT improves throughput by

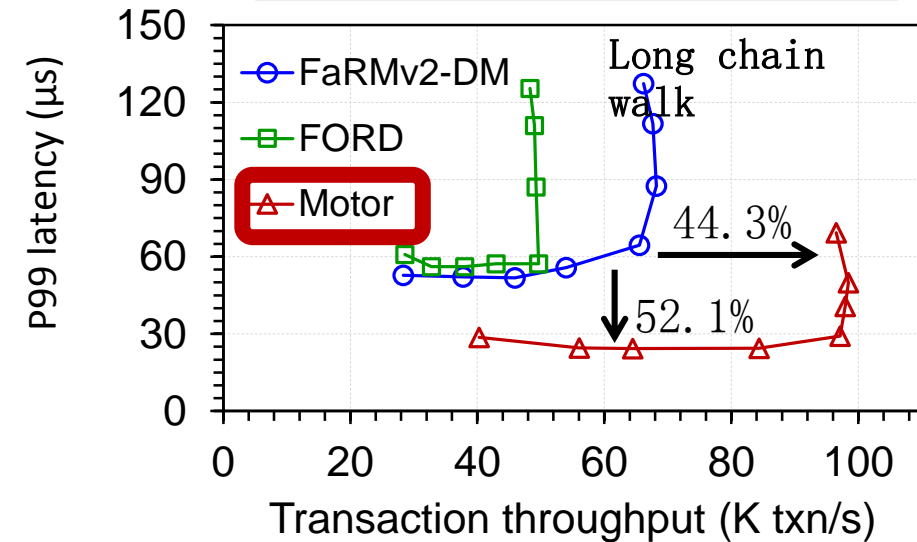$1.7 - 2.4x$ over O2N

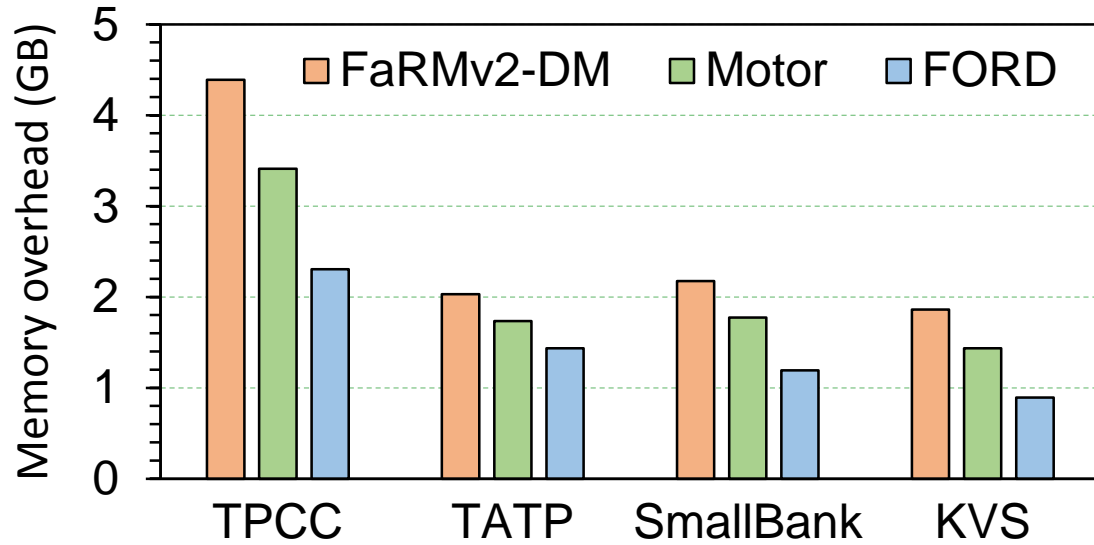$1.3 - 1.6x$ over

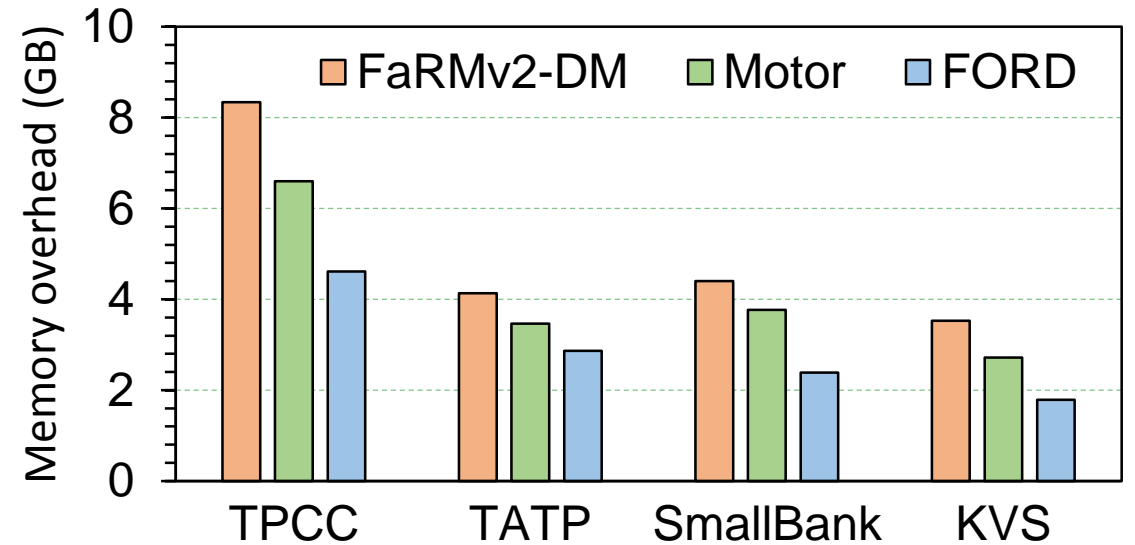# End-to-End Performance



TATP (read-intensive)

TPCC (write-intensive)

# Memory Overhead



Memory consumption on workload scale1*

Memory consumption on workload scale2*

Motor's full-delta value storage design

- VS. FORD (single-versioning): 1.45x, not 4x, on TPCC with 4 versions
- VS. FaRMv2 (multi-versioning): save 14.6%~22.8% of memory space
- Stable advantages on larger scale of workloads

* Workload scale1: TPCC 24 warehouses; TATP 2M subscribers; SmallBank 10M accounts; KVS 10M objects
* Workload scale2: TPCC 48 warehouses; TATP 4M subscribers; SmallBank 20M accounts; KVS 20M objects

# Conclusion

➢ Distributed transaction is a key pillar for disaggregated memory

➢ Limitations of existing systems

  • Single-versioning: limited concurrency, high logging overhead
  • Multi-versioning: inefficient linked chain, incompatible txn protocol

➢ *Motor:* a holistic multi-versioning design

  • Memory pool: *consecutive version tuple*
  • Compute pool: *one-sided RDMA-based MVCC*

➢ **Benefits**

**High Throughput**    **Low Latency**    **Low Memory Overhead**

https://github.com/minghust/motor

# Thank you! Q&A