

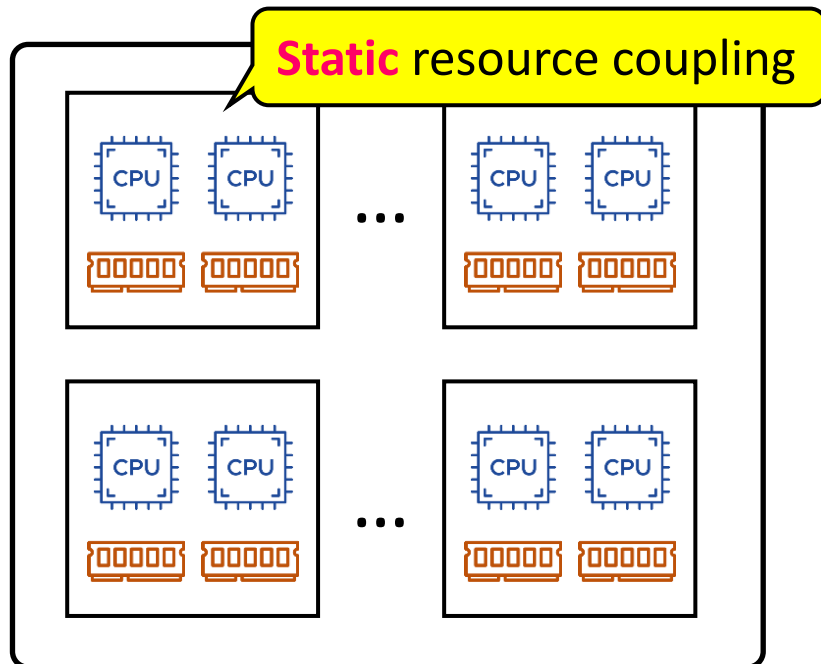
# Motor: Enabling Multi-Versioning for Distributed Transactions on Disaggregated Memory

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

*Huazhong University of Science and Technology, China*

# Insufficient Memory Utilization in Cloud

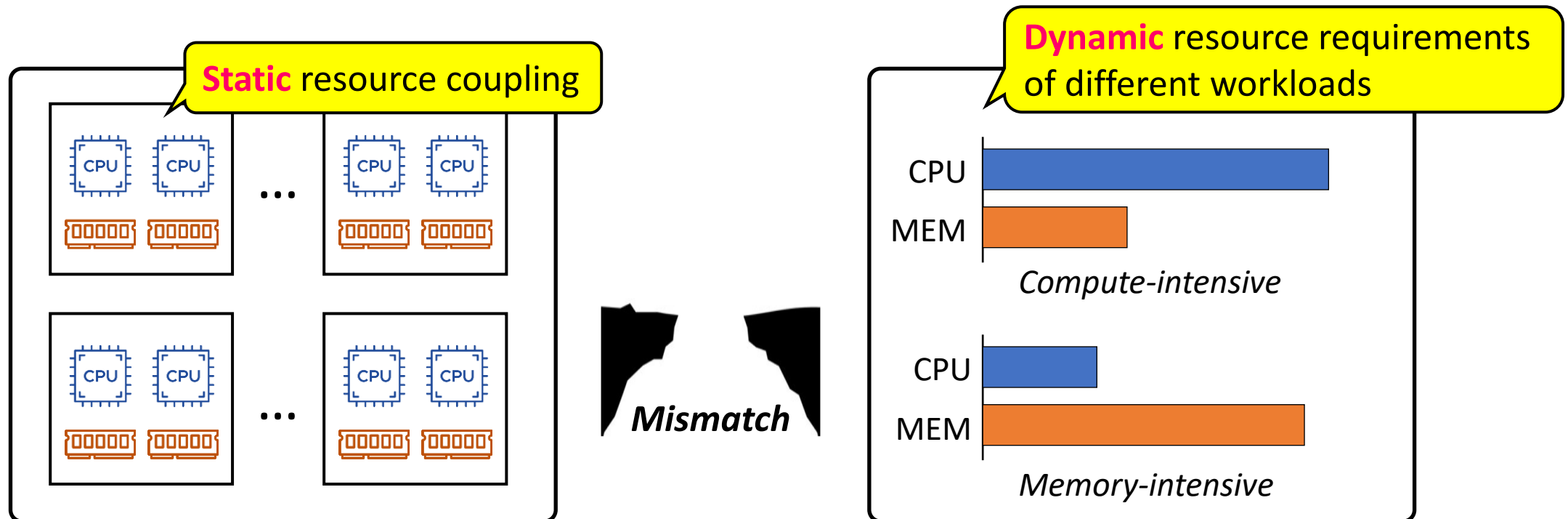
- About 20%~60% [1-4]
- One major reason: monolithic server



- [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>

# Insufficient Memory Utilization in Cloud

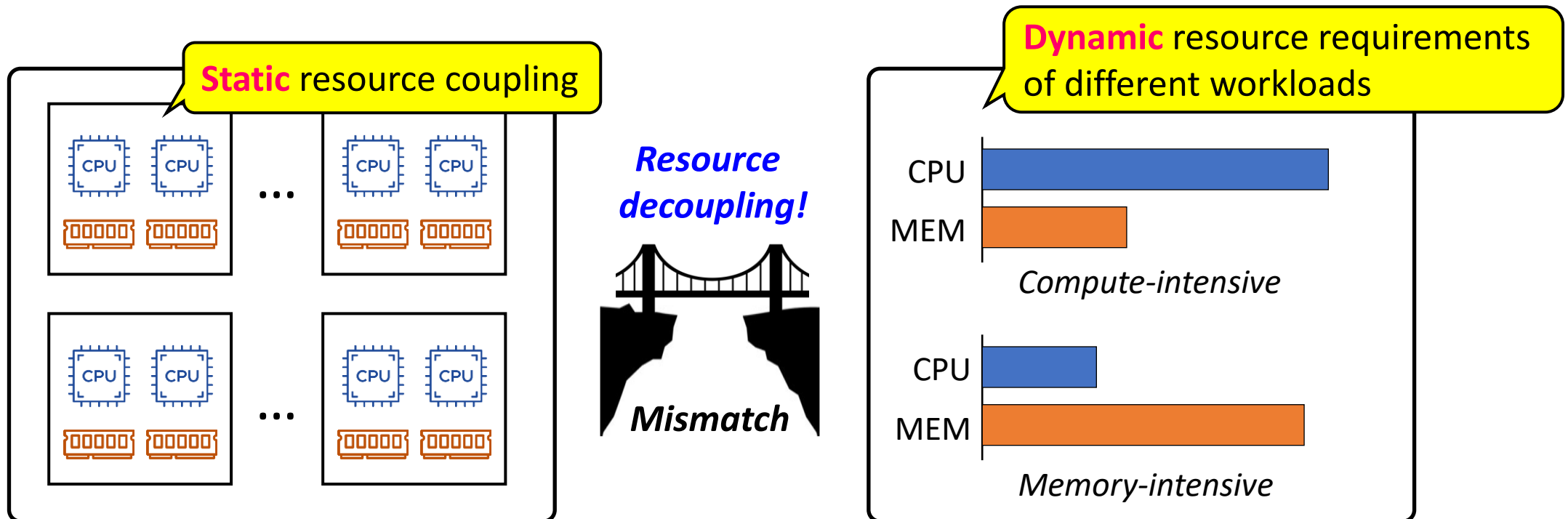
- About 20%~60% [1-4]
- One major reason: monolithic server



- [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>

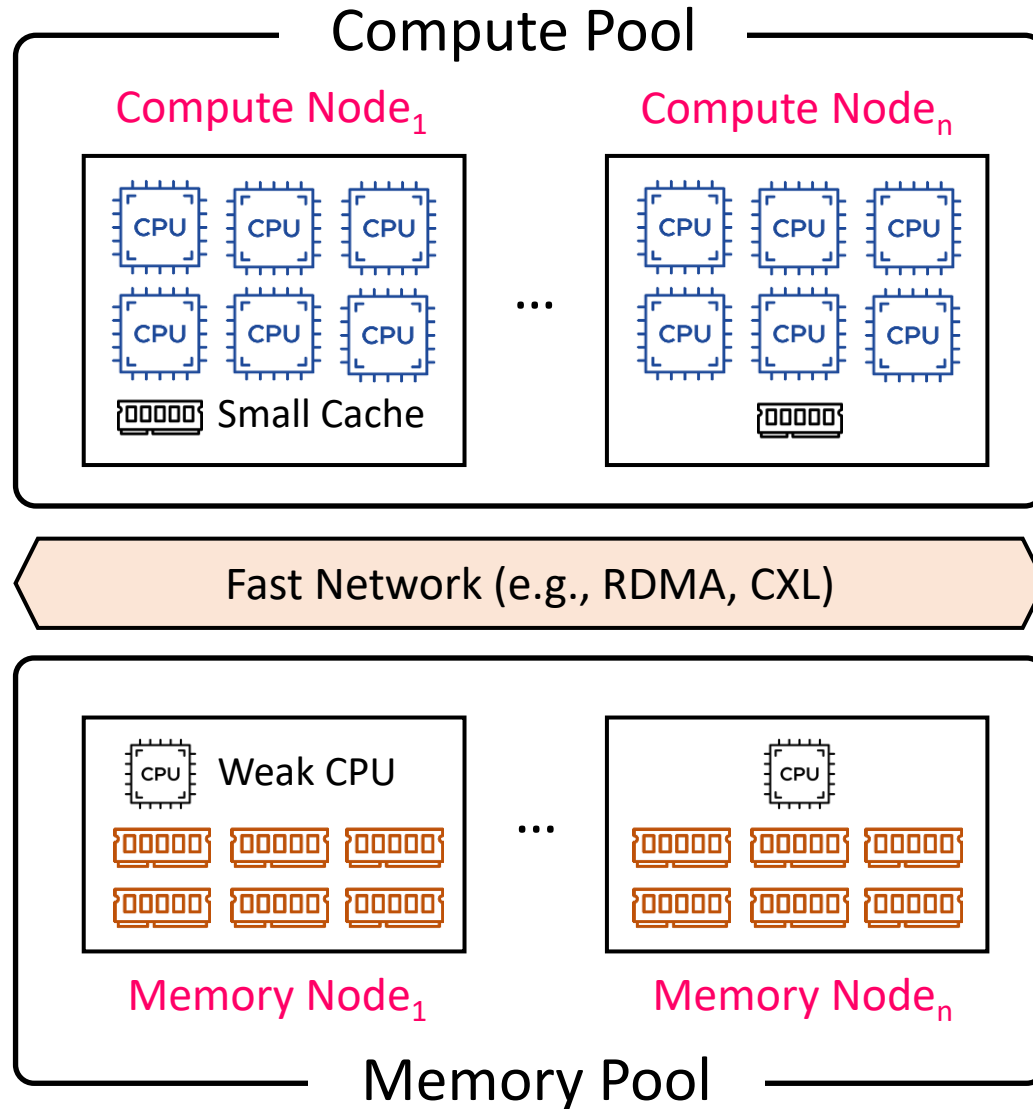
# Insufficient Memory Utilization in Cloud

- About 20%~60% [1-4]
- One major reason: monolithic server

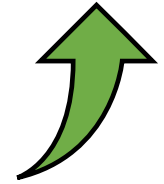


- [1] MemTrade@SIGMETRICS'23, Borg@EuroSys'20, LegoOS@OSDI'18  
[2] Google Production Cluster Trace. <https://github.com/google/cluster-data>  
[3] Alibaba Production Cluster Trace. <https://github.com/alibaba/clusterdata>  
[4] Snowflake Dataset. <https://github.com/resource-disaggregation/snowset>

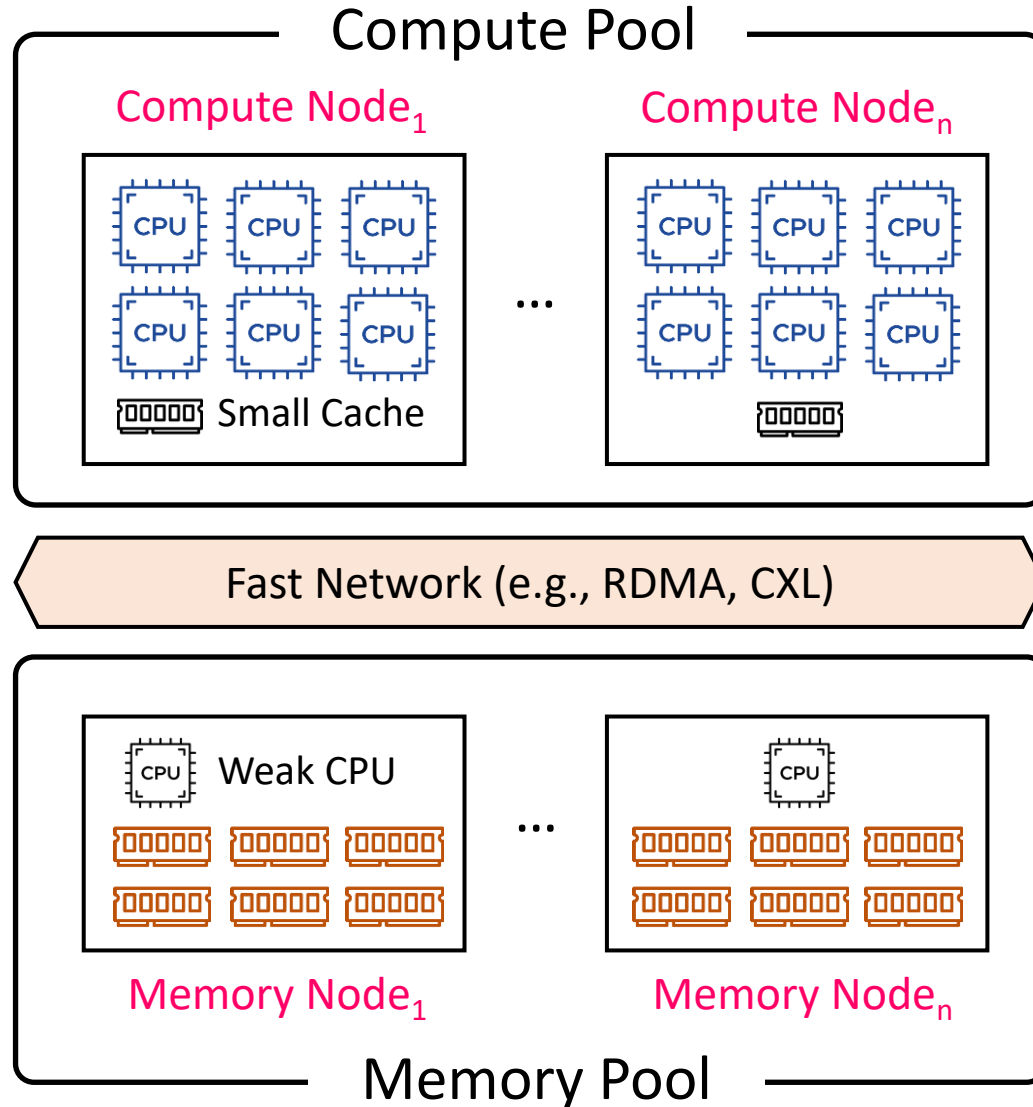
# Disaggregated Memory (DM)



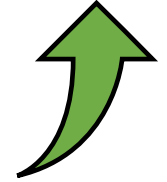
*Resource utilization*  
*Elasticity*



# Disaggregated Memory (DM)

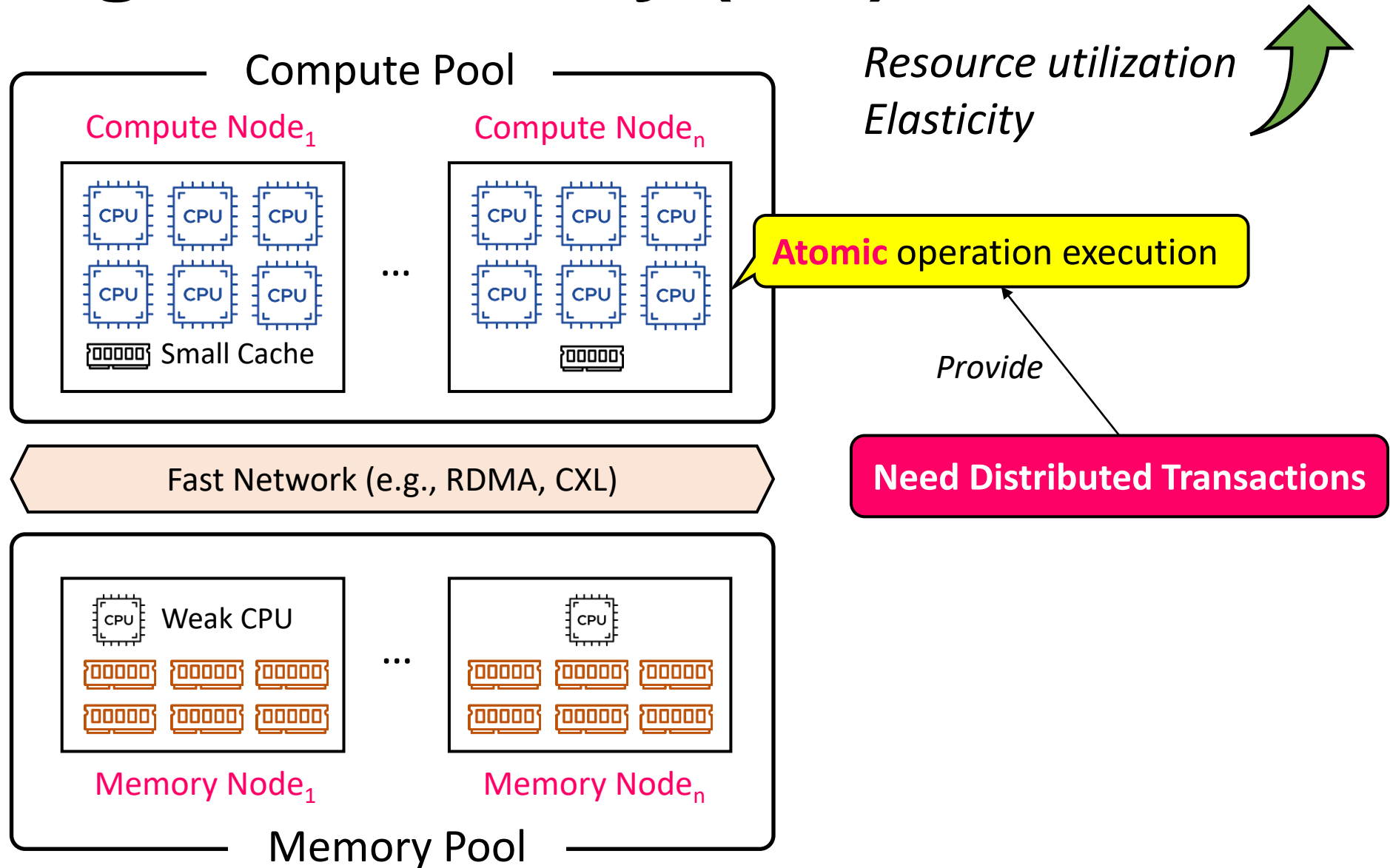


*Resource utilization*  
*Elasticity*

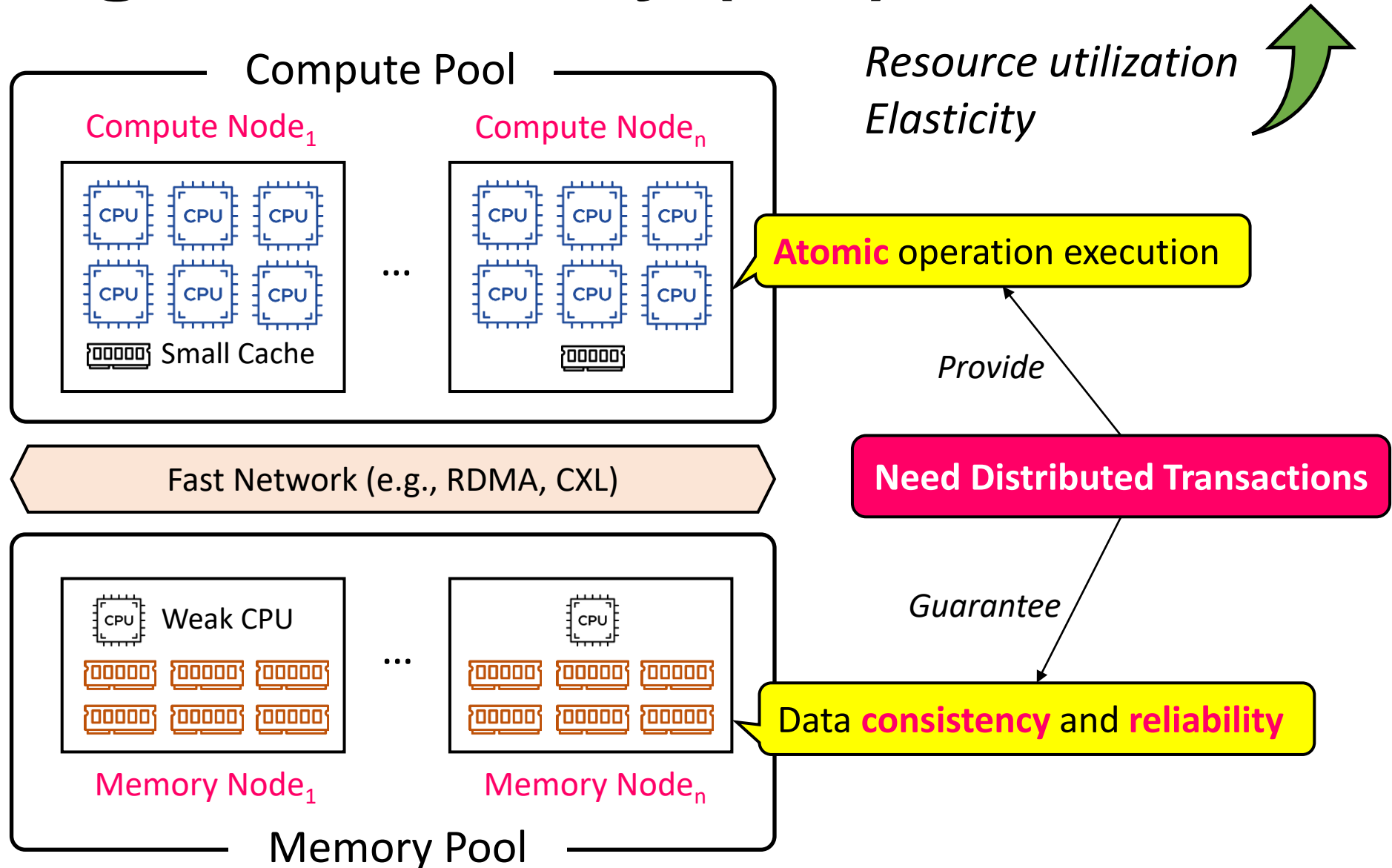


**Need Distributed Transactions**

# Disaggregated Memory (DM)

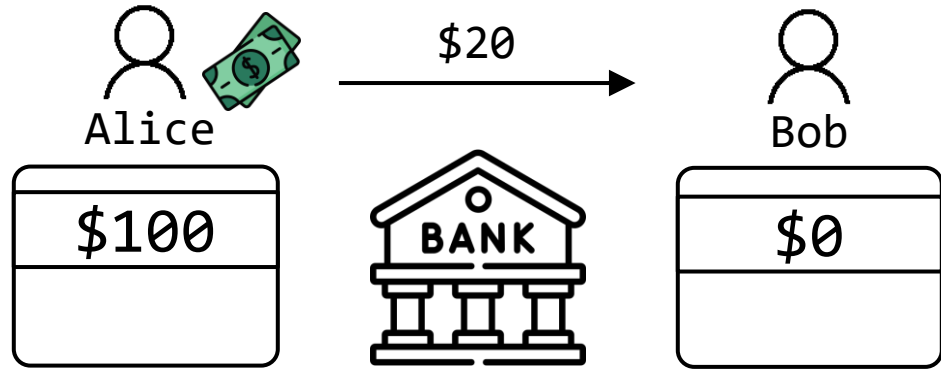


# Disaggregated Memory (DM)

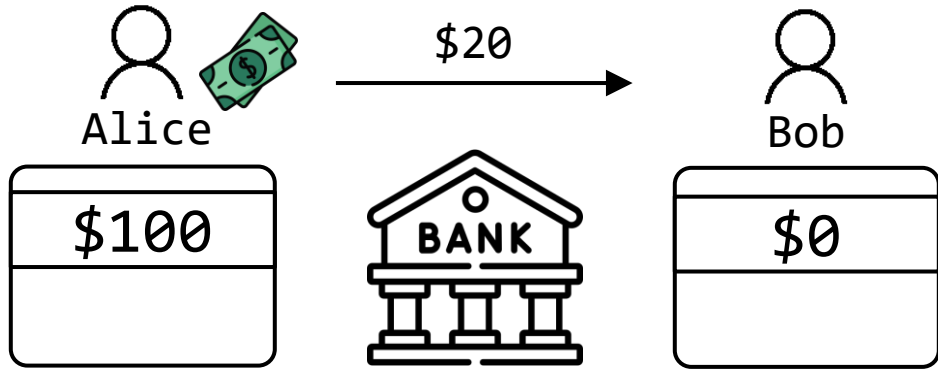




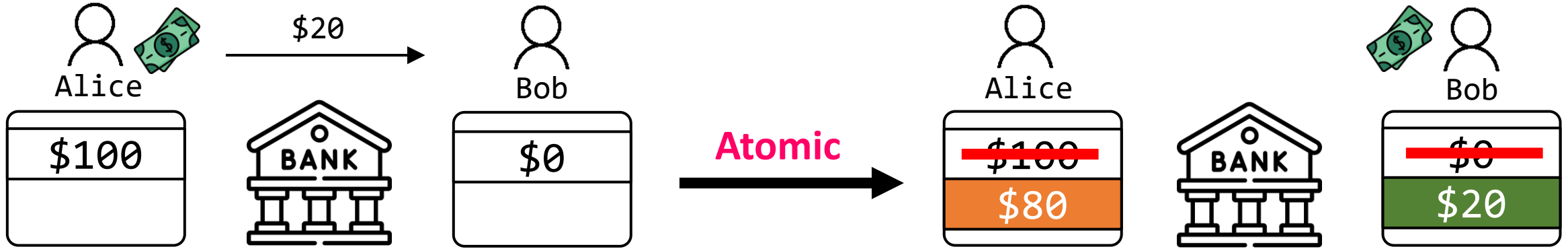
# Transaction



# Transaction



# Transaction



*Txn begin*

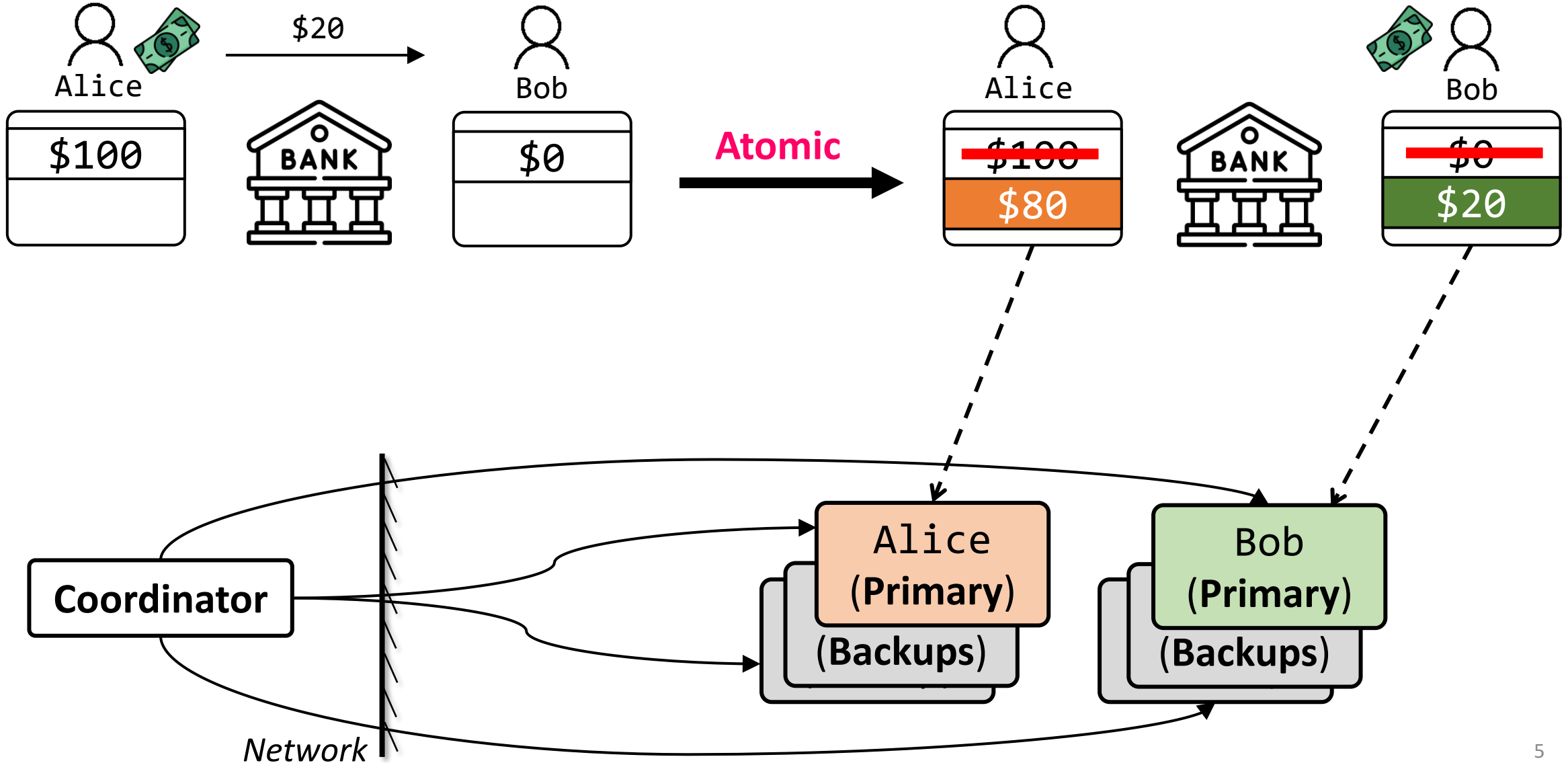
Alice: \$100 → \$80

Bob: \$0 → \$20

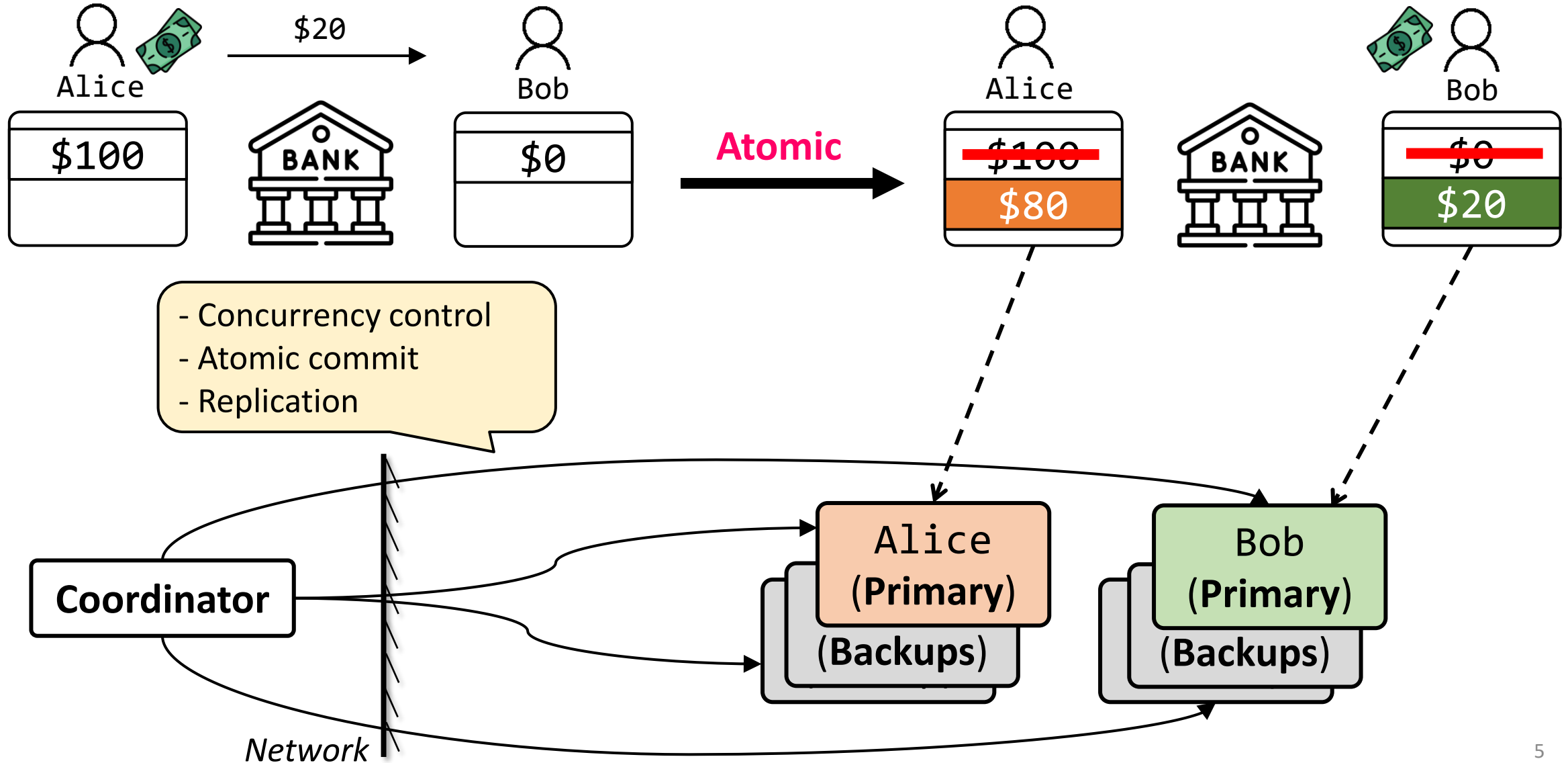
*Txn end*

*Transaction*

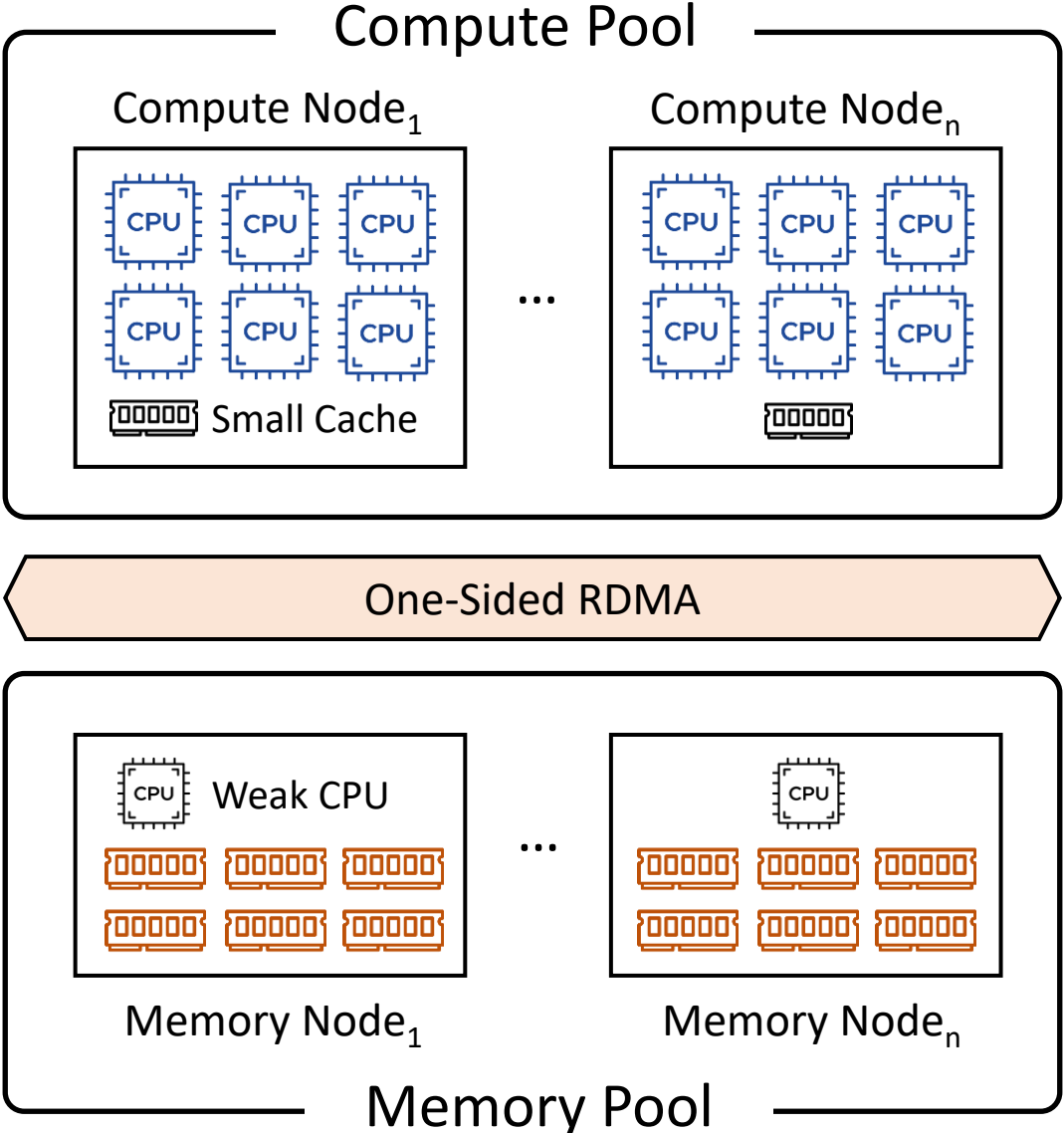
# Distributed Transaction



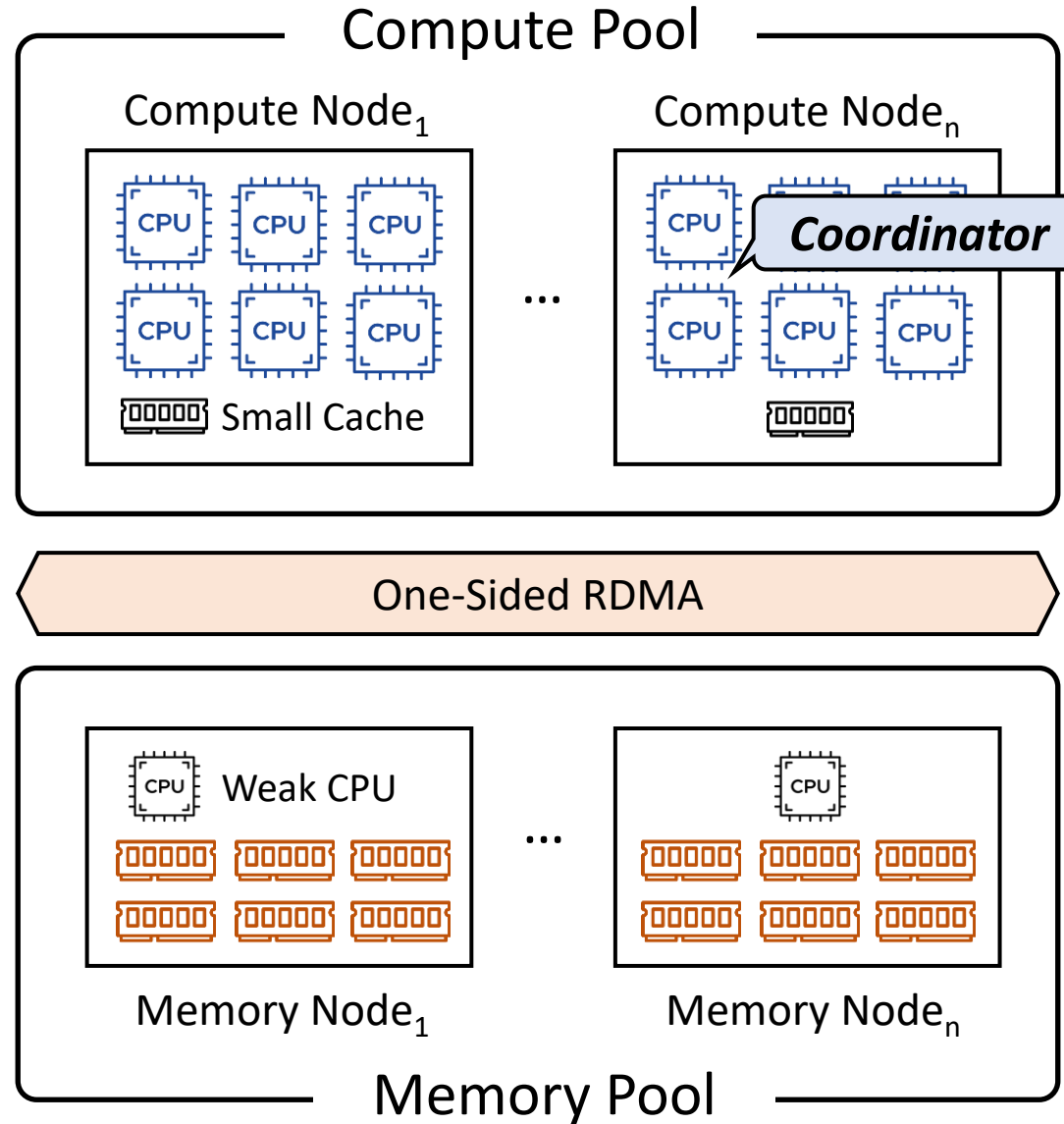
# Distributed Transaction



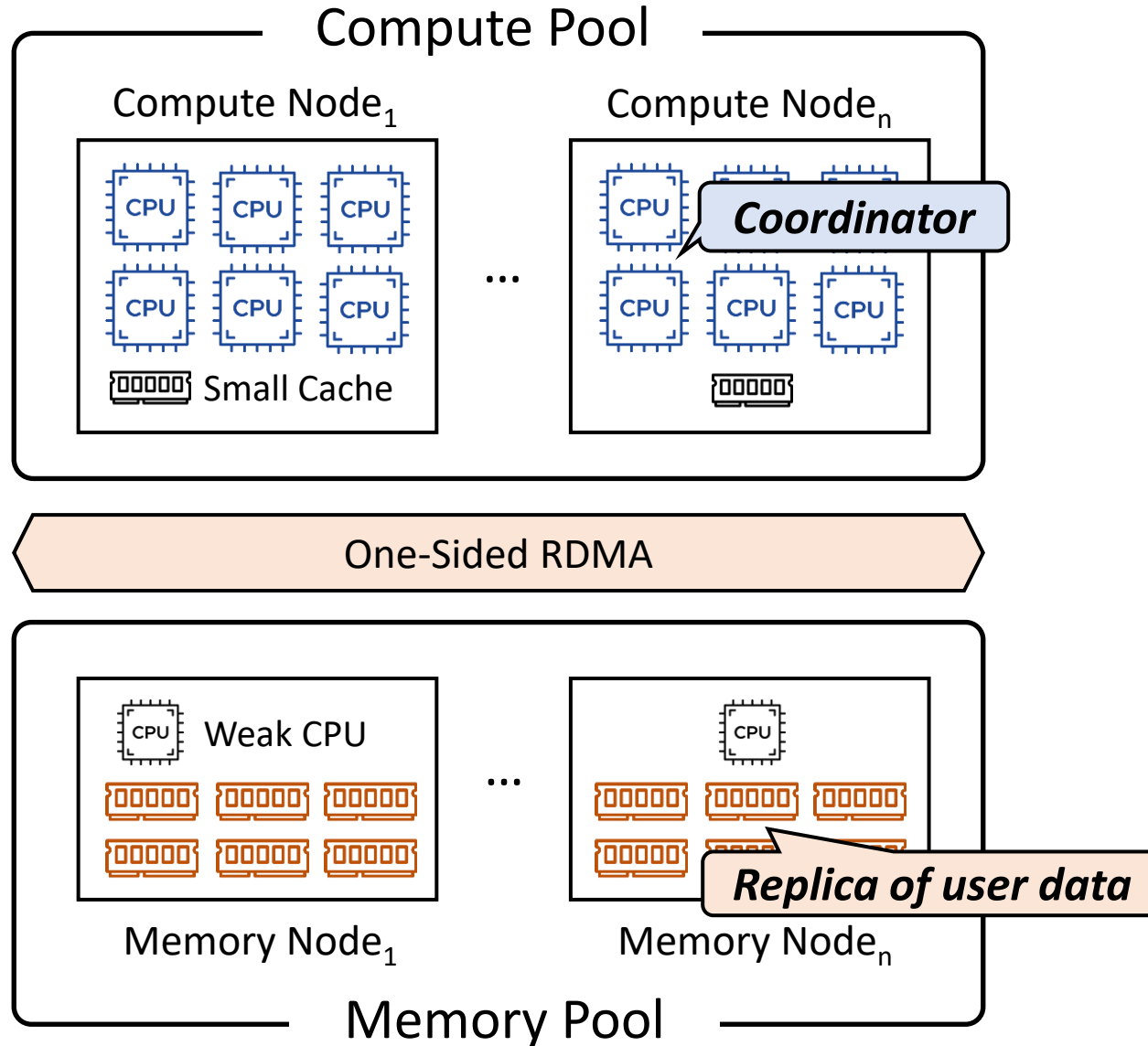
# Distributed Transaction on DM



# Distributed Transaction on DM

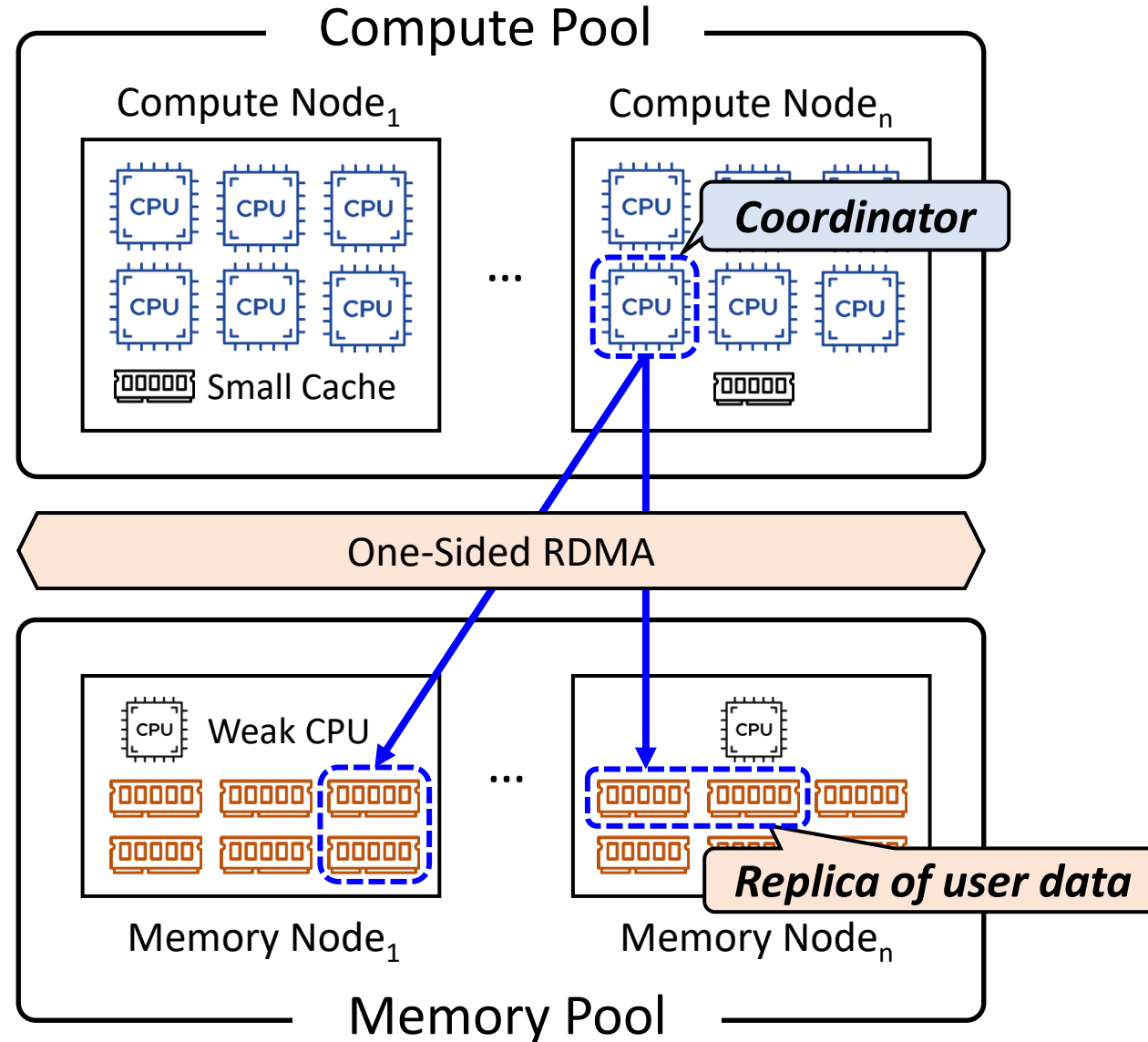


# Distributed Transaction on DM

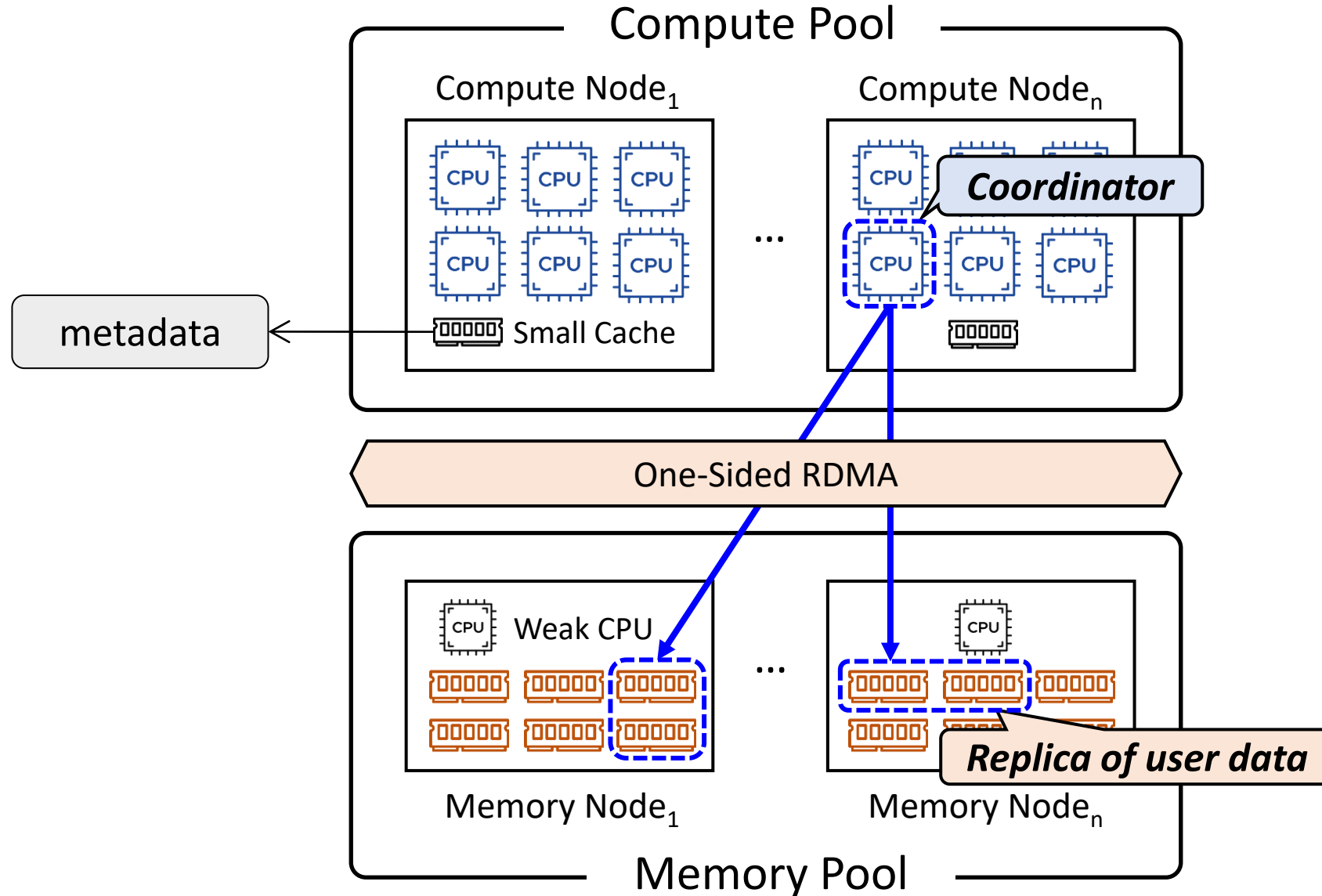




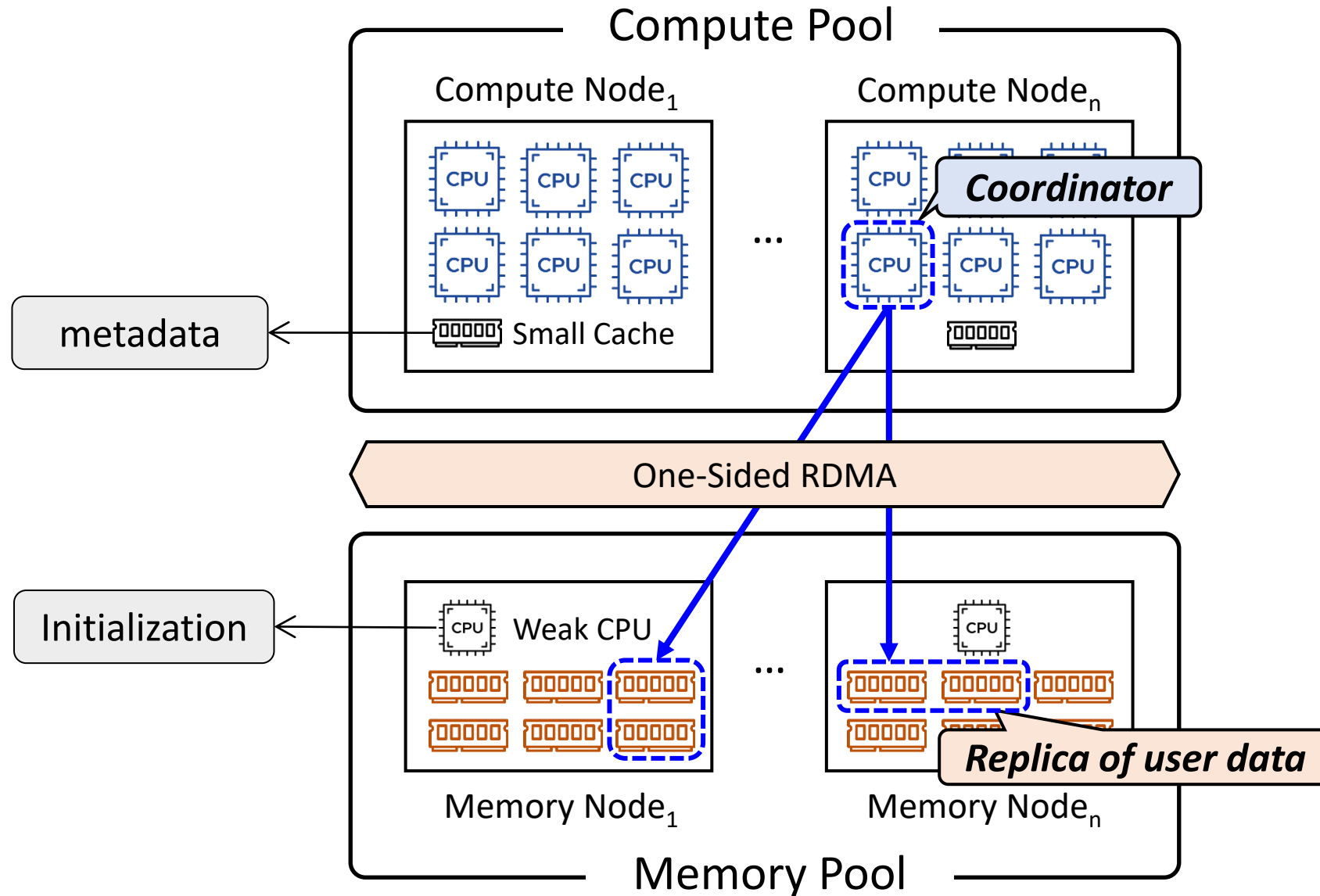
# Distributed Transaction on DM



# Distributed Transaction on DM



# Distributed Transaction on DM



# Existing Studies

- Single-versioning distributed transaction system for DM<sup>[1]</sup>

# Existing Studies

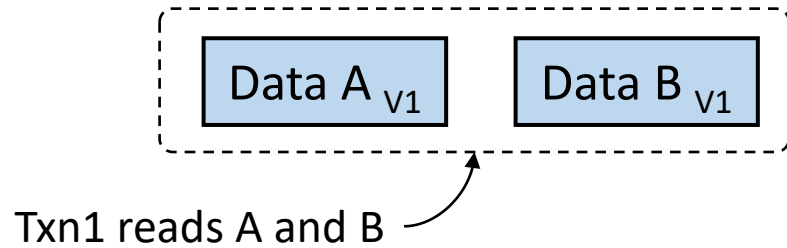
➤ Single-versioning distributed transaction system for DM<sup>[1]</sup>

☹ Limited concurrency

# Existing Studies

- Single-versioning distributed transaction system for DM<sup>[1]</sup>

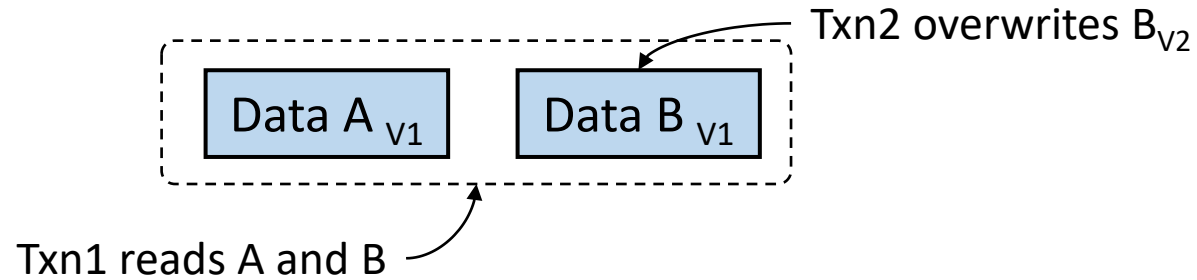
☹ Limited concurrency



# Existing Studies

- Single-versioning distributed transaction system for DM<sup>[1]</sup>

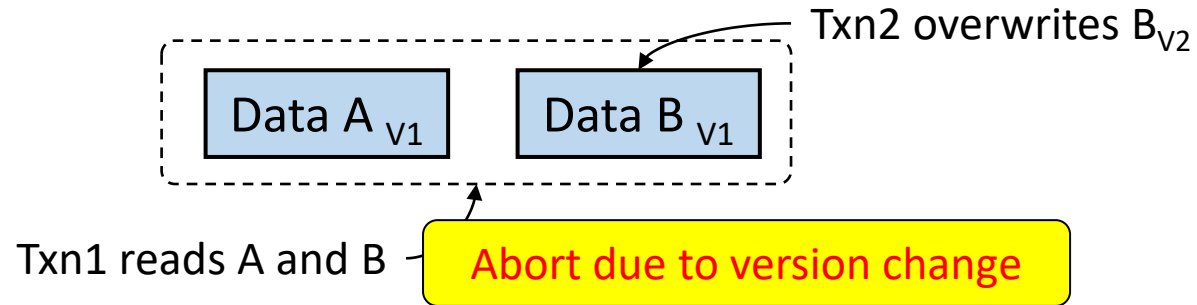
☹ Limited concurrency



# Existing Studies

- Single-versioning distributed transaction system for DM<sup>[1]</sup>

☹ Limited concurrency



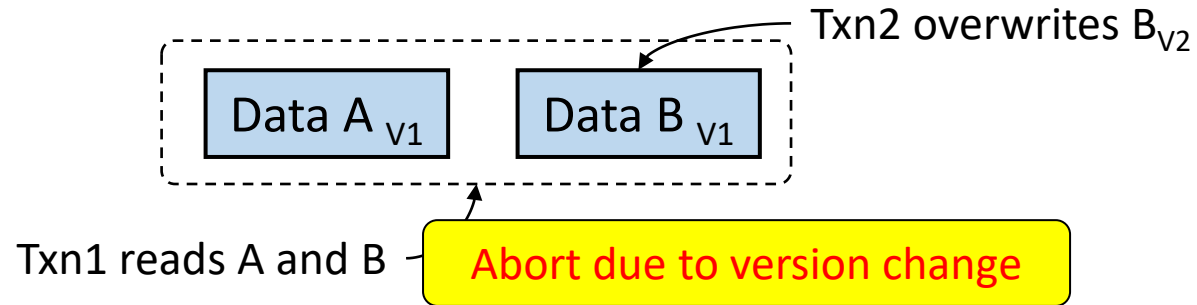


# Existing Studies

- Single-versioning distributed transaction system for DM<sup>[1]</sup>

☹ Limited concurrency

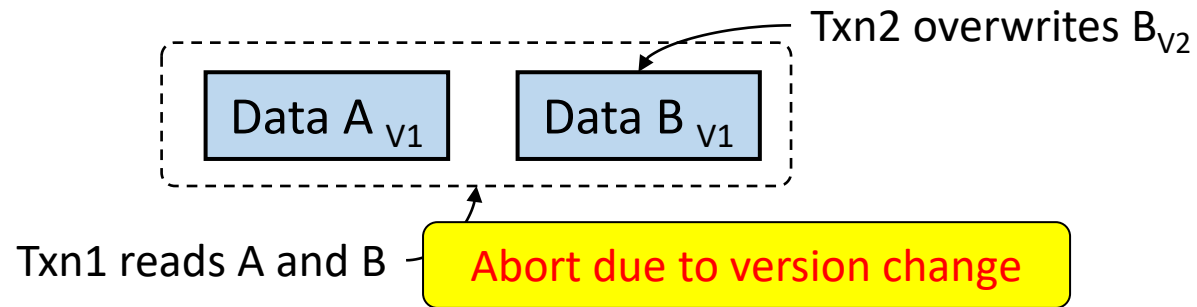
☹ High logging overhead



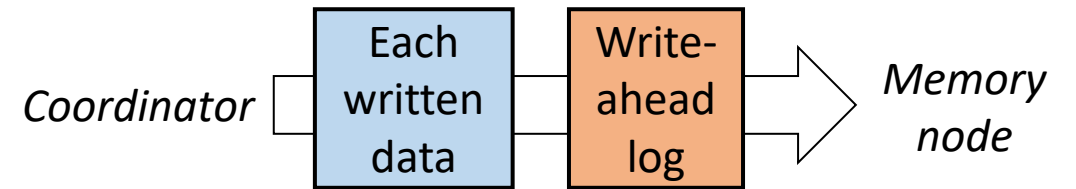
# Existing Studies

- Single-versioning distributed transaction system for DM<sup>[1]</sup>

☹ Limited concurrency



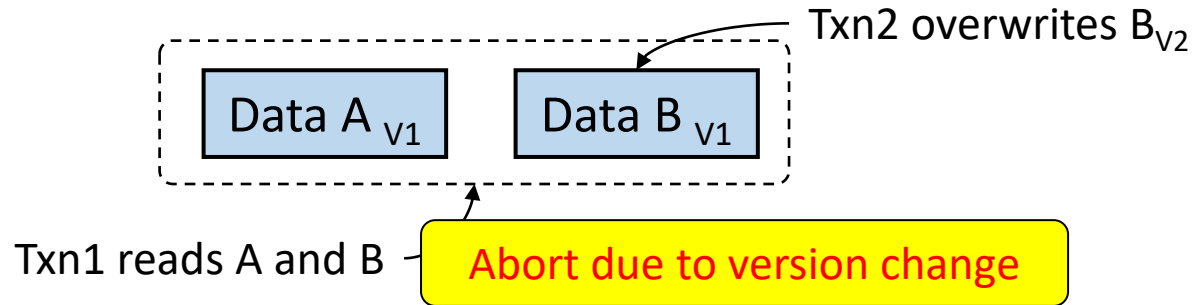
☹ High logging overhead



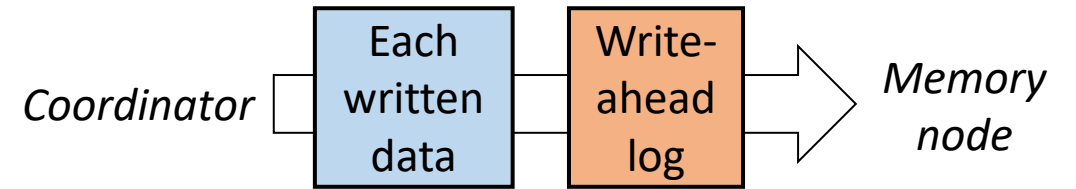
# Existing Studies

- Single-versioning distributed transaction system for DM<sup>[1]</sup>

☹ Limited concurrency



☹ High logging overhead

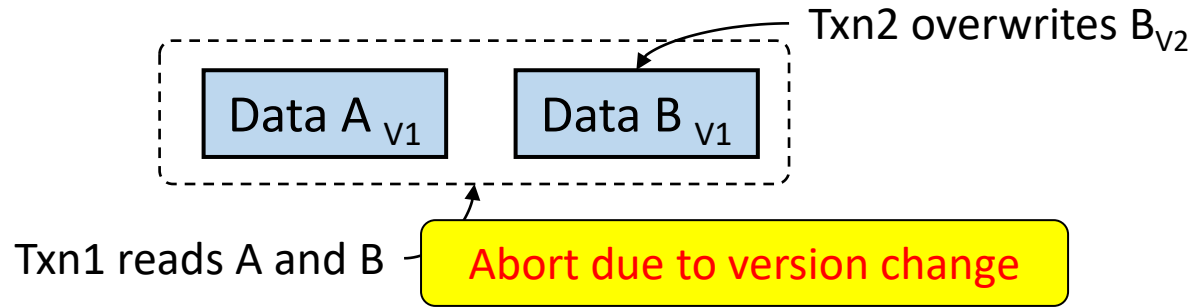


- Multi-versioning helps address limitations of single-versioning

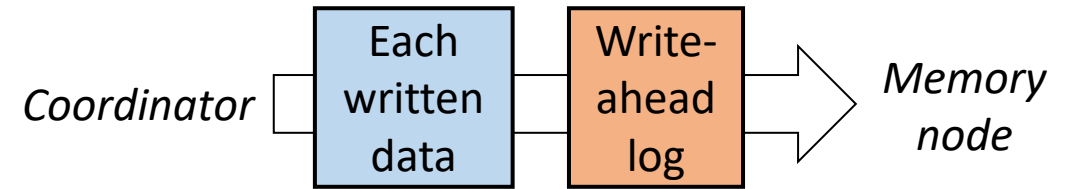
# Existing Studies

- Single-versioning distributed transaction system for DM<sup>[1]</sup>

☹️ Limited concurrency

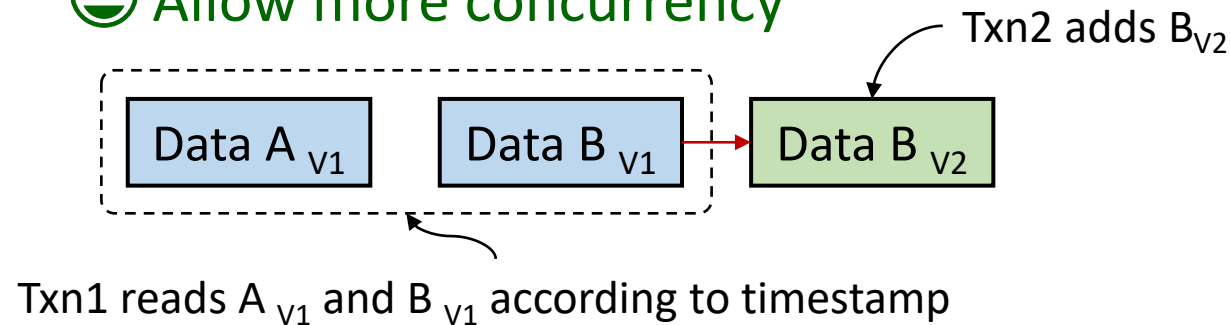


☹️ High logging overhead



- Multi-versioning helps address limitations of single-versioning

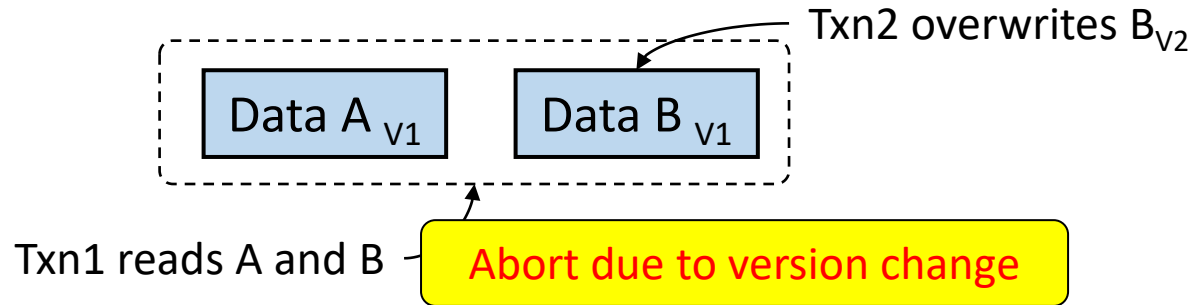
😊 Allow more concurrency



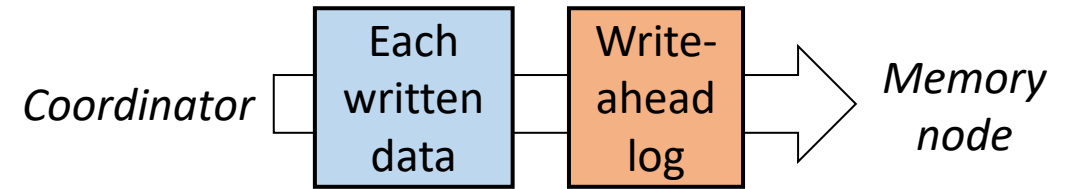
# Existing Studies

- Single-versioning distributed transaction system for DM<sup>[1]</sup>

☹️ Limited concurrency

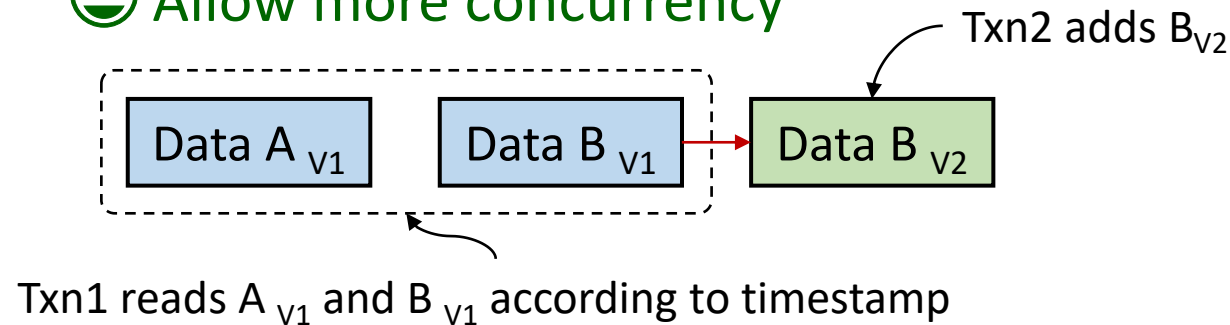


☹️ High logging overhead

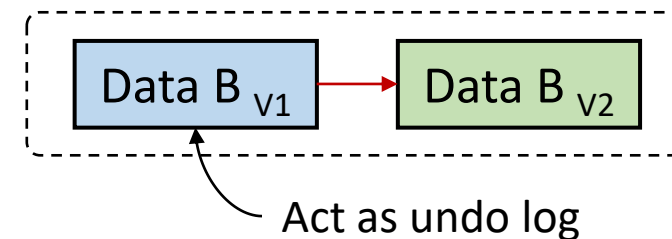


- Multi-versioning helps address limitations of single-versioning

😊 Allow more concurrency



😊 Reduce logs



# Does Multi-Versioning Work on DM?

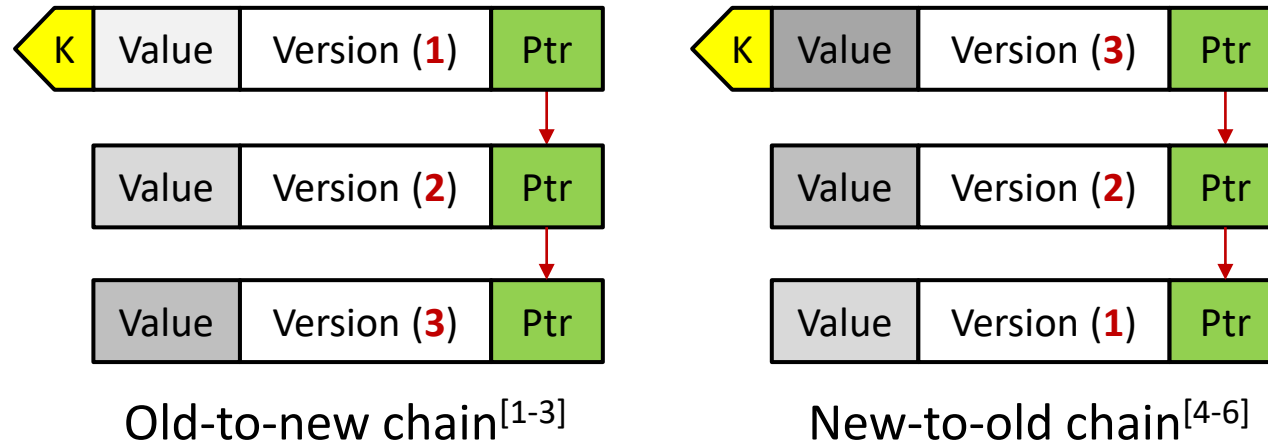
- Existing systems are based on monolithic servers

# Does Multi-Versioning Work on DM?

➤ Existing systems are based on monolithic servers

☹️ Inefficient linked version chain

- Works in monolithic servers but does not fit DM



[1] DST@NSDI'21 [2] Hekaton@SIGMOD'13 [3] Aurogon@FAST'22

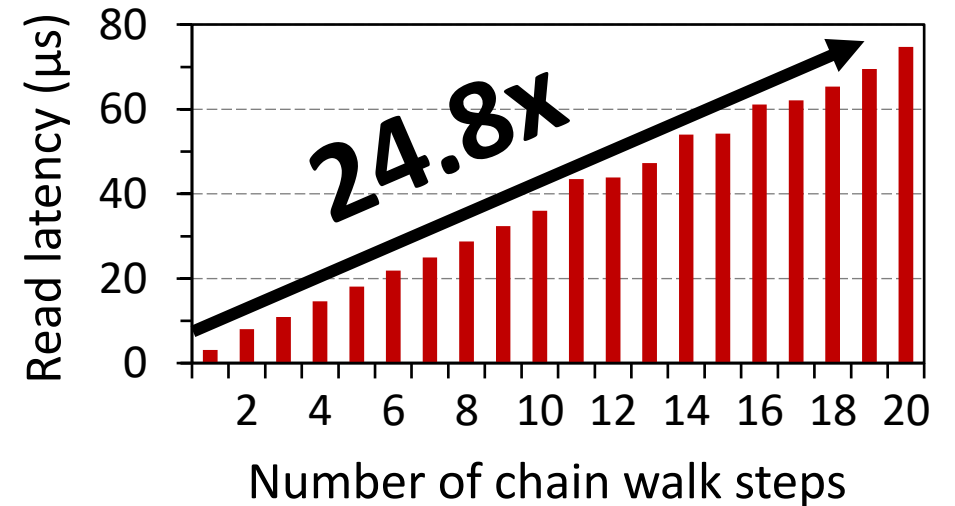
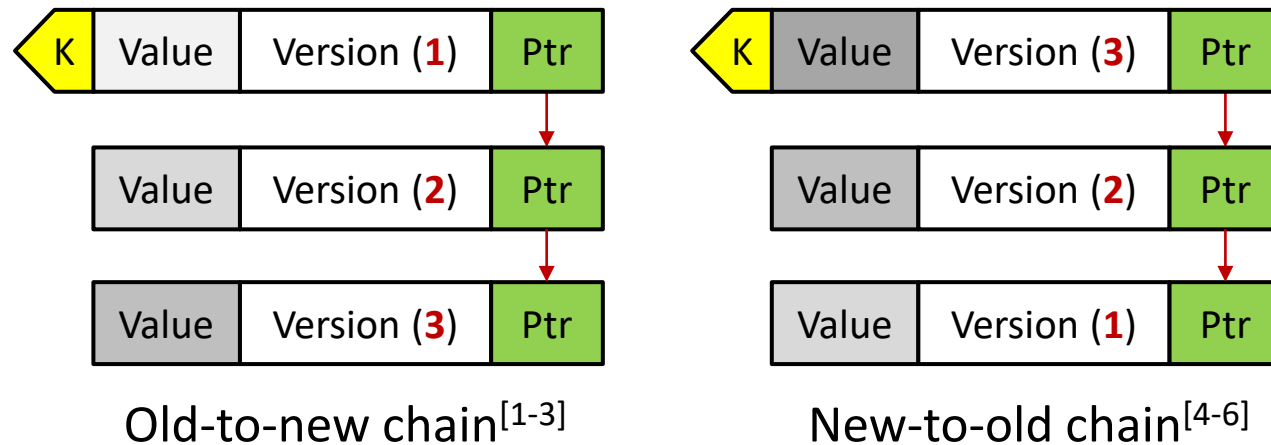
[4] FaRMv2@SIGMOD'19 [5] Neumann et al.@SIGMOD'15 [6] NAM-DB@VLDB'17

# Does Multi-Versioning Work on DM?

➤ Existing systems are based on monolithic servers

☹️ **Inefficient linked version chain**

- Works in monolithic servers but does not fit DM



[1] DST@NSDI'21 [2] Hekaton@SIGMOD'13 [3] Aurogon@FAST'22

[4] FaRMv2@SIGMOD'19 [5] Neumann et al.@SIGMOD'15 [6] NAM-DB@VLDB'17

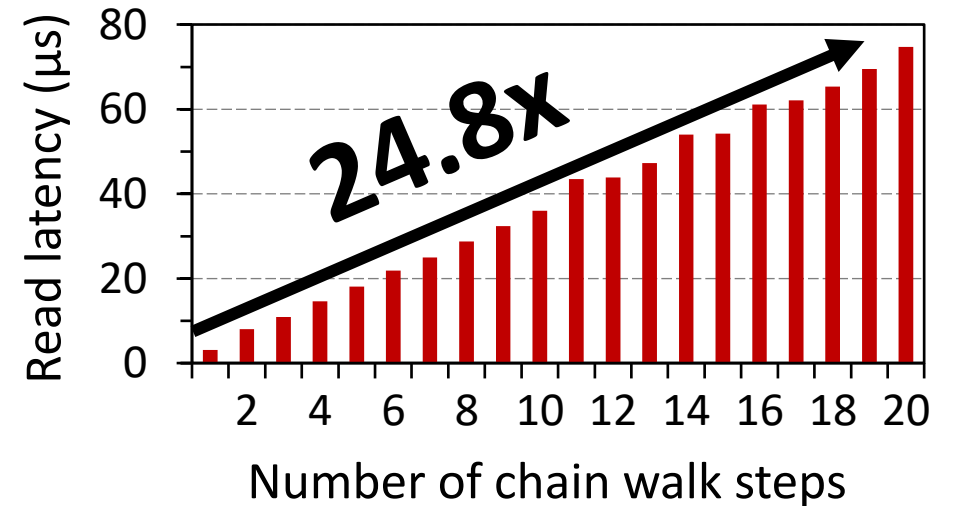
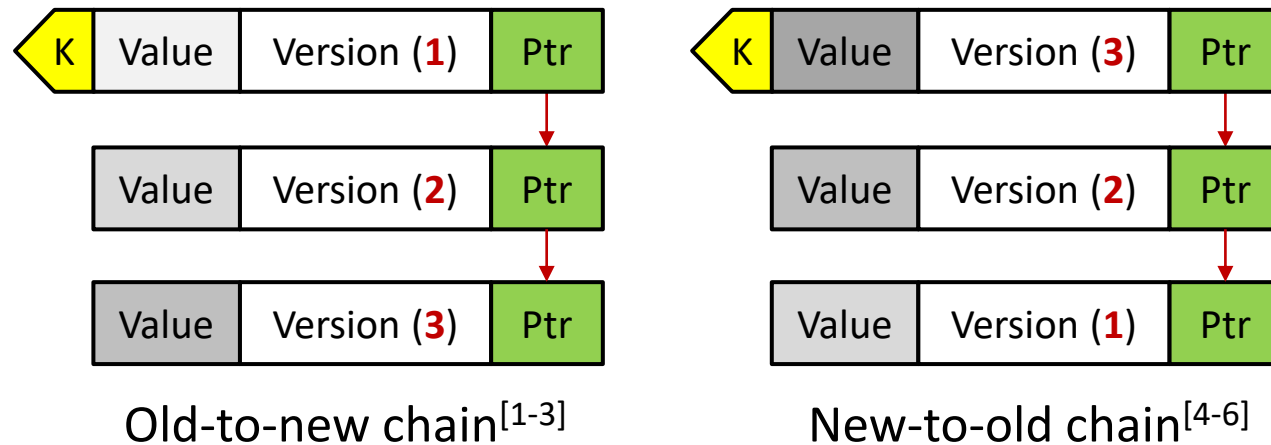


# Does Multi-Versioning Work on DM?

➤ Existing systems are based on monolithic servers

## ☹️ Inefficient linked version chain

- Works in monolithic servers but does not fit DM



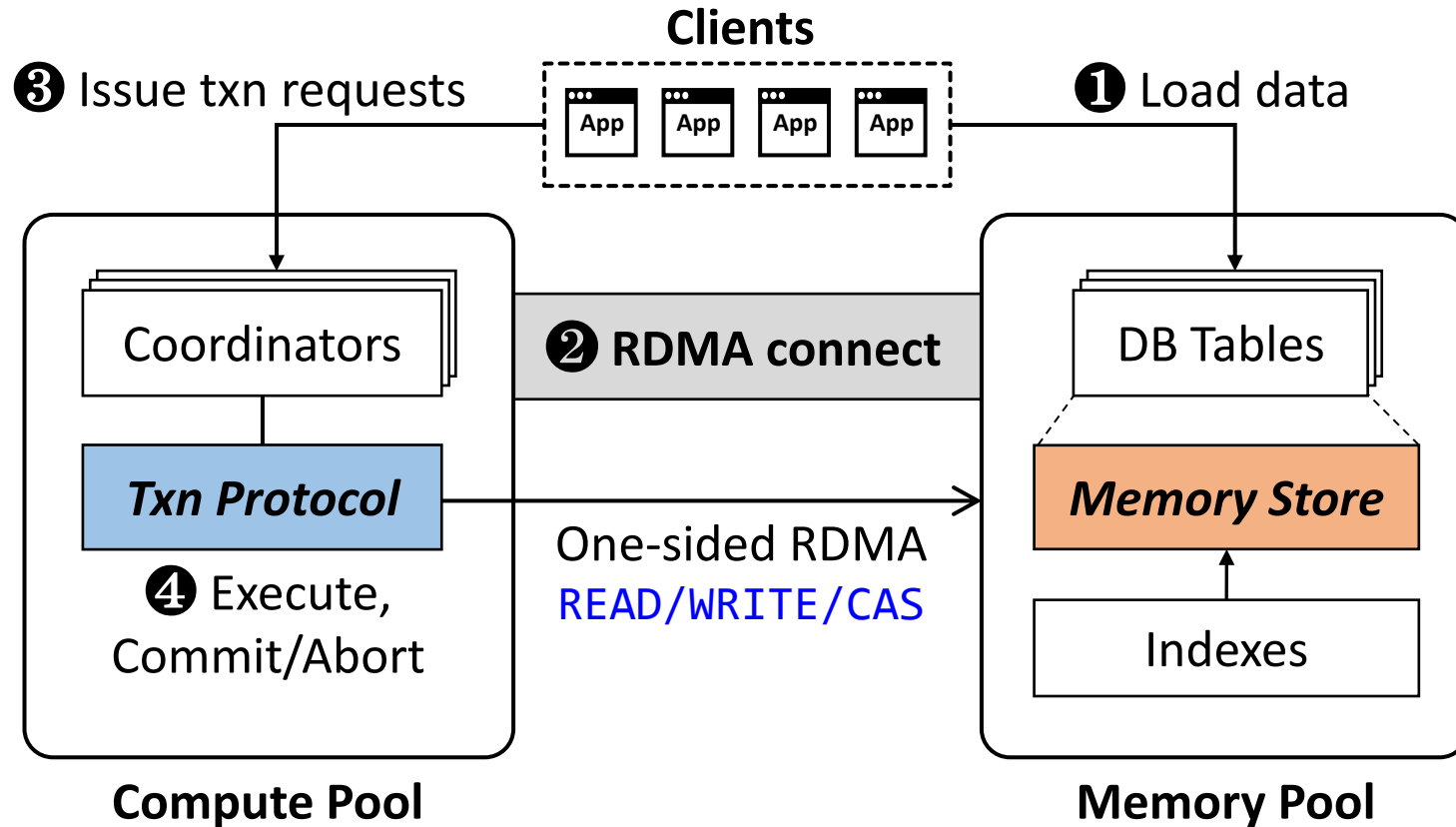
## ☹️ Incompatible transaction protocol

- Frequently consumes CPU of each data node<sup>[1,4,5]</sup>
- Memory node stores data but only has weak CPU

[1] DST@NSDI'21 [2] Hekaton@SIGMOD'13 [3] Aurogon@FAST'22

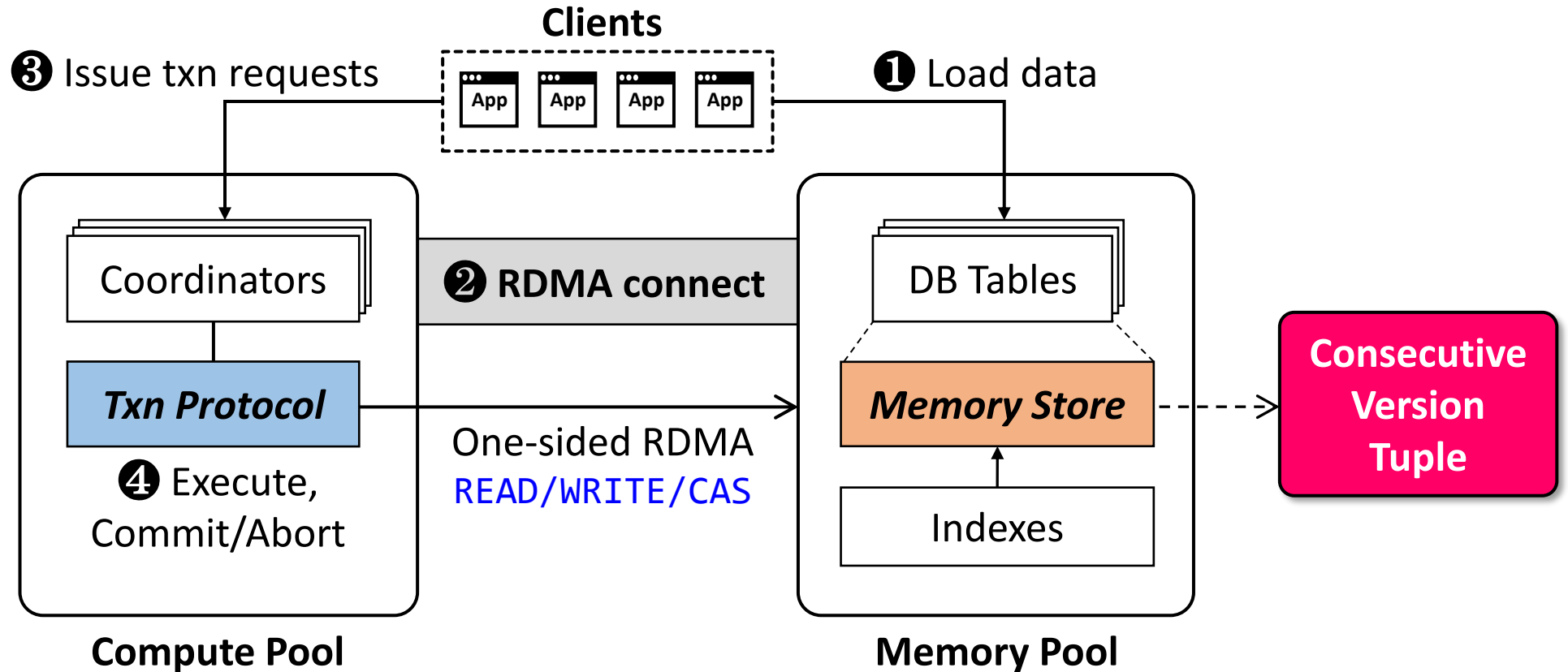
[4] FaRMv2@SIGMOD'19 [5] Neumann et al.@SIGMOD'15 [6] NAM-DB@VLDB'17

# Motor: Enabling Multi-Versioning for DM



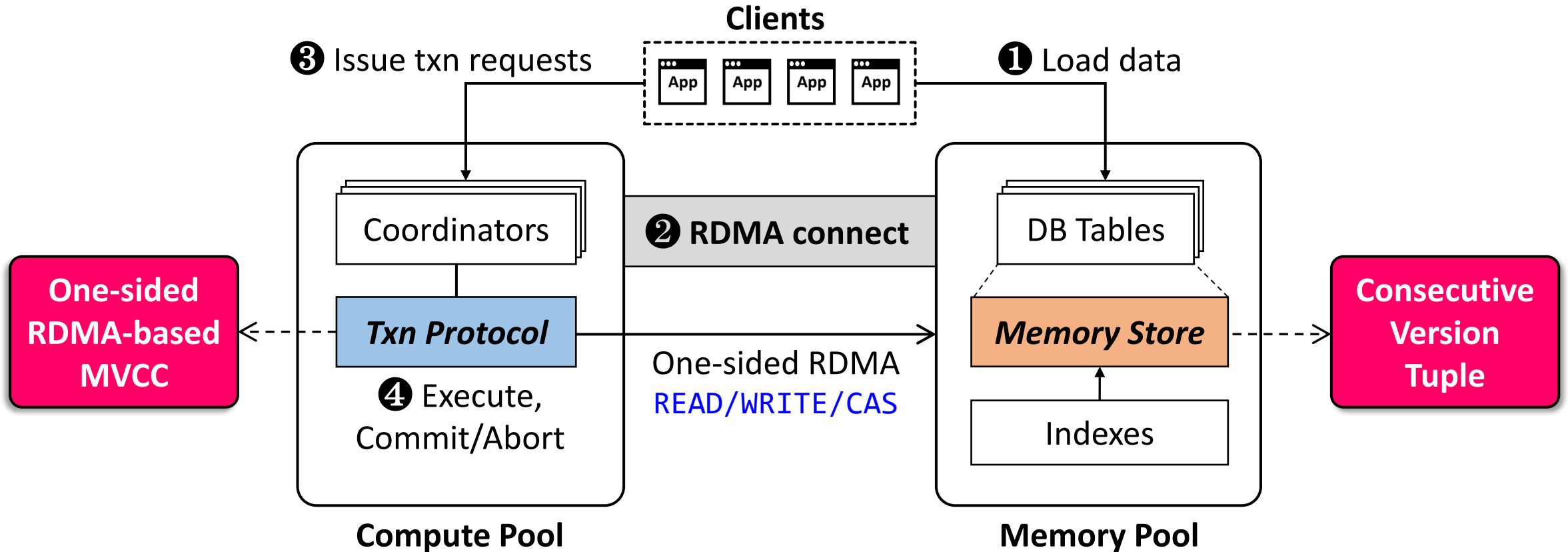
System Overview

# Motor: Enabling Multi-Versioning for DM



System Overview

# Motor: Enabling Multi-Versioning for DM



System Overview

# Consecutive Version Tuple

CVT Region

A consecutive version tuple (CVT)

Header <sub>1</sub>	Vcell <sub>1</sub>	...	Vcell <sub>n-1</sub>	Vcell <sub>n</sub>
Header <sub>2</sub>	Vcell <sub>1</sub>	...	Vcell <sub>n-1</sub>	Vcell <sub>n</sub>
...	...	...	...	...
Header <sub>n</sub>	Vcell <sub>1</sub>	...	Vcell <sub>n-1</sub>	Vcell <sub>n</sub>



# Consecutive Version Tuple

CVT Region

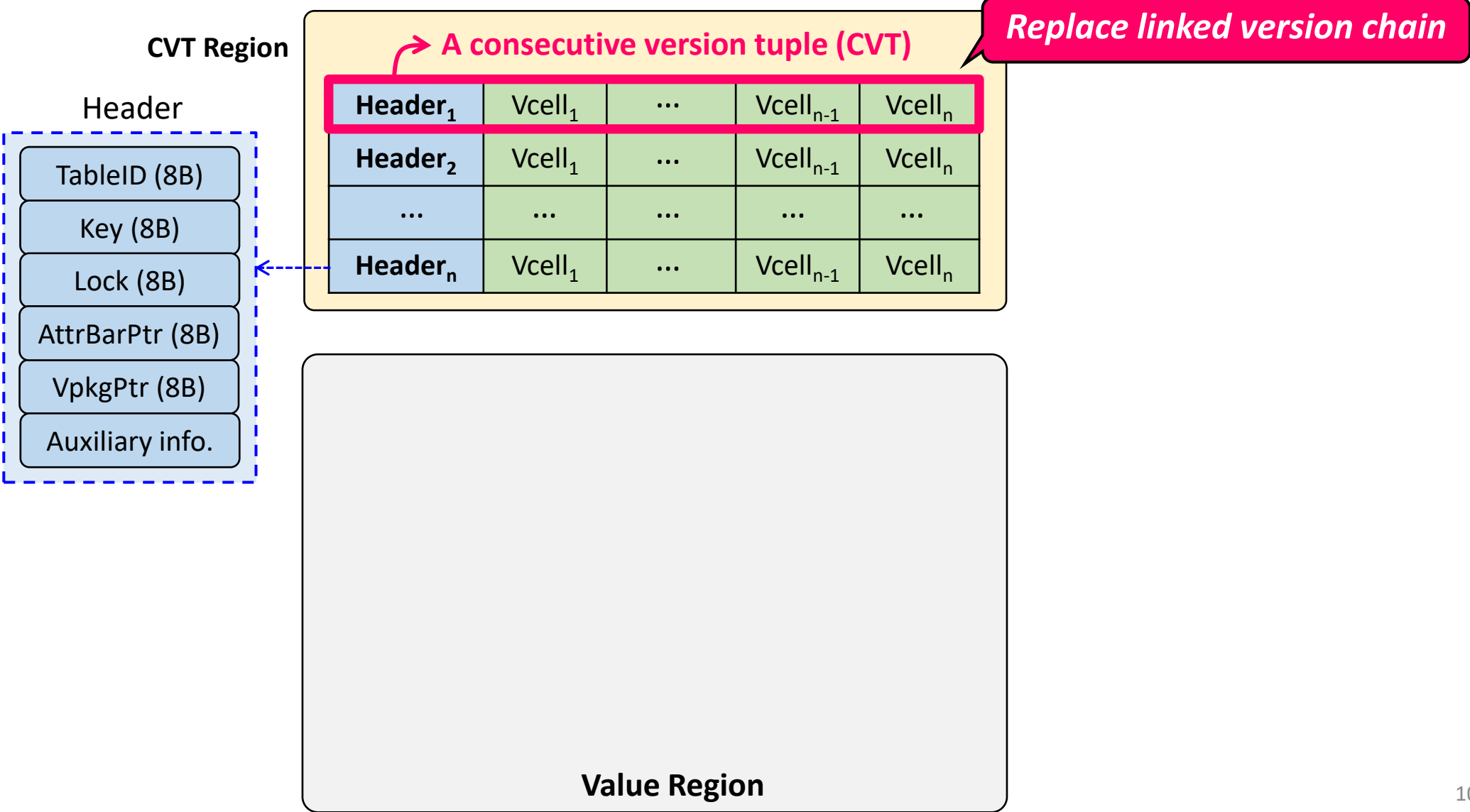
A consecutive version tuple (CVT)

Header <sub>1</sub>	Vcell <sub>1</sub>	...	Vcell <sub>n-1</sub>	Vcell <sub>n</sub>
Header <sub>2</sub>	Vcell <sub>1</sub>	...	Vcell <sub>n-1</sub>	Vcell <sub>n</sub>
...	...	...	...	...
Header <sub>n</sub>	Vcell <sub>1</sub>	...	Vcell <sub>n-1</sub>	Vcell <sub>n</sub>

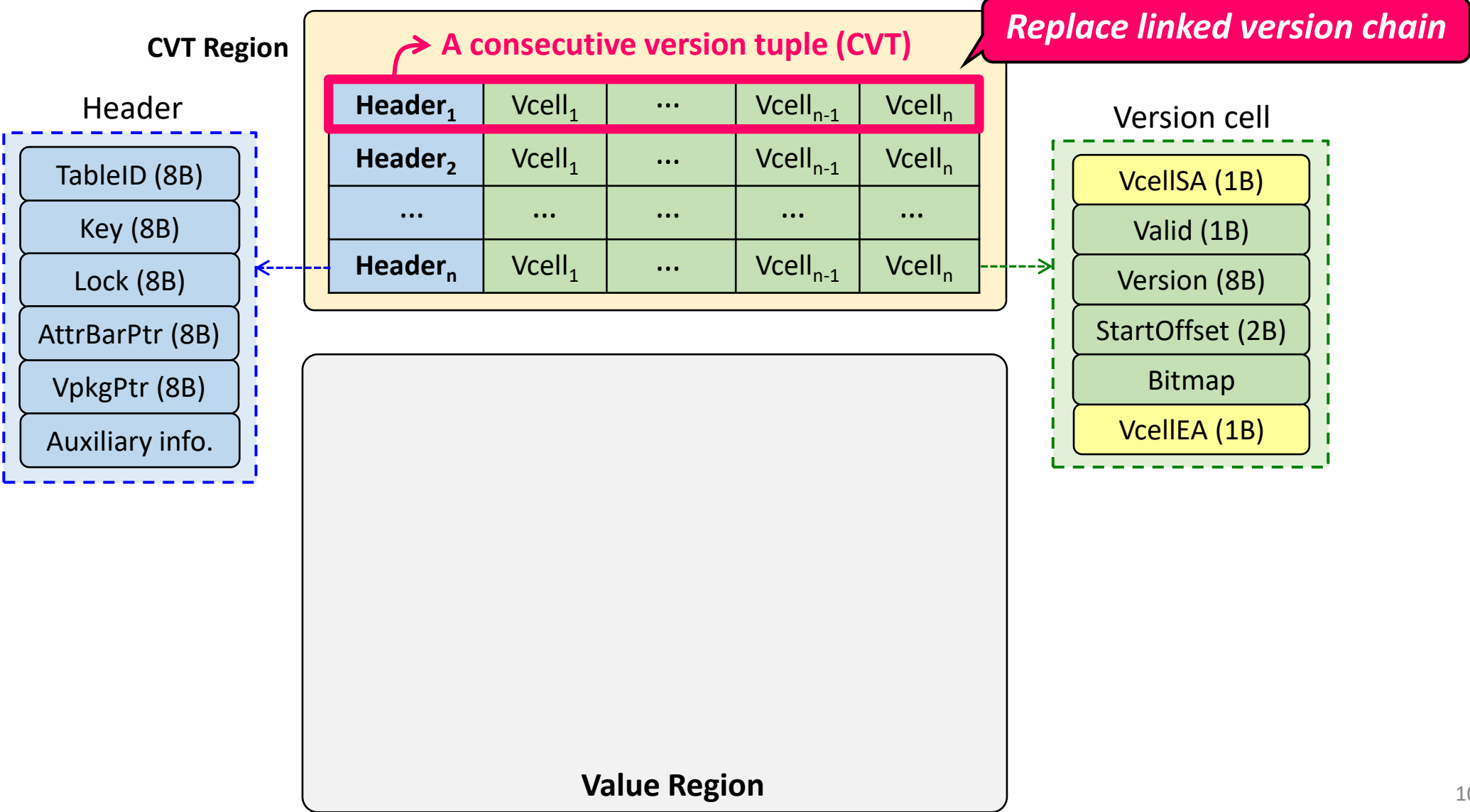
Replace linked version chain



# Consecutive Version Tuple

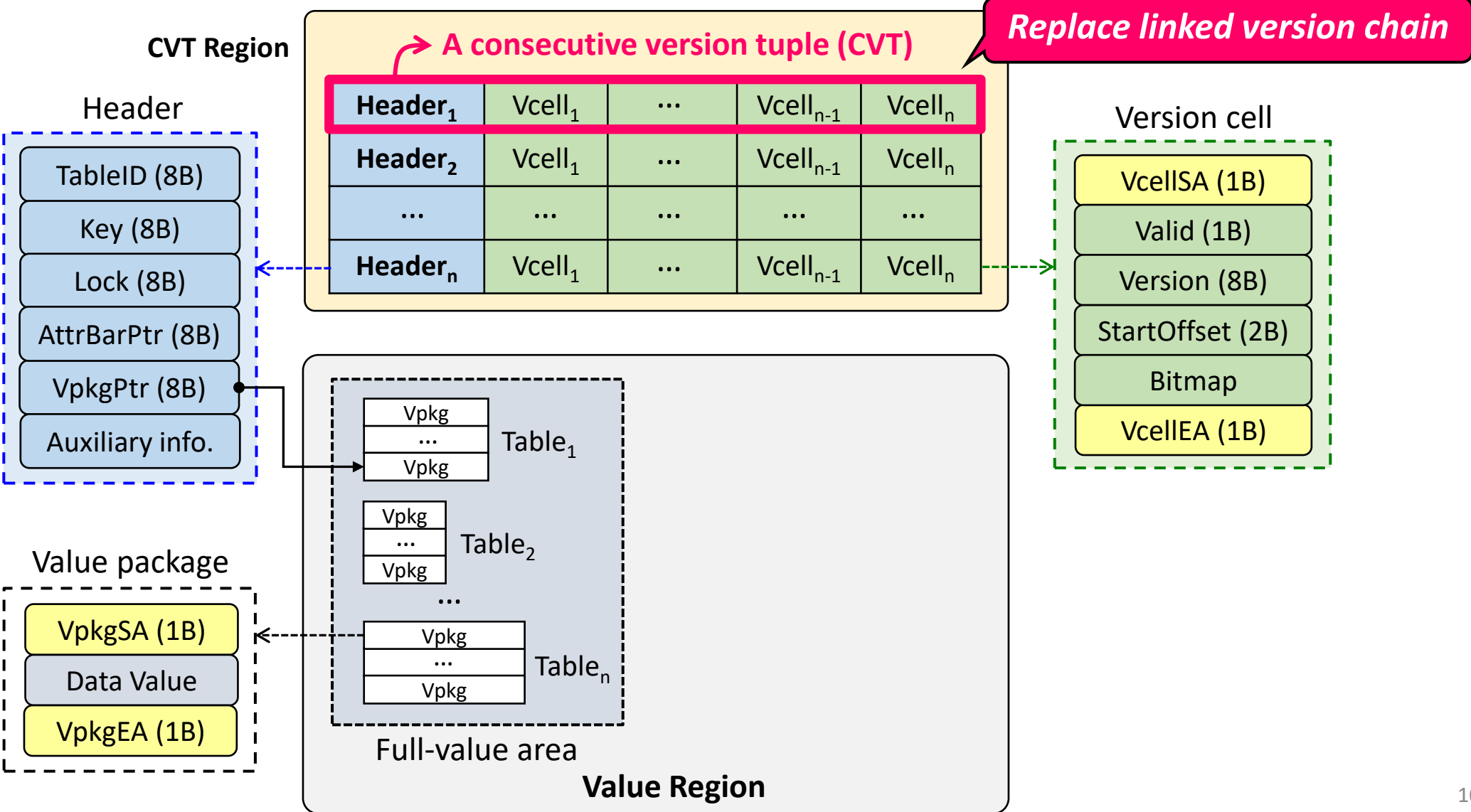


# Consecutive Version Tuple

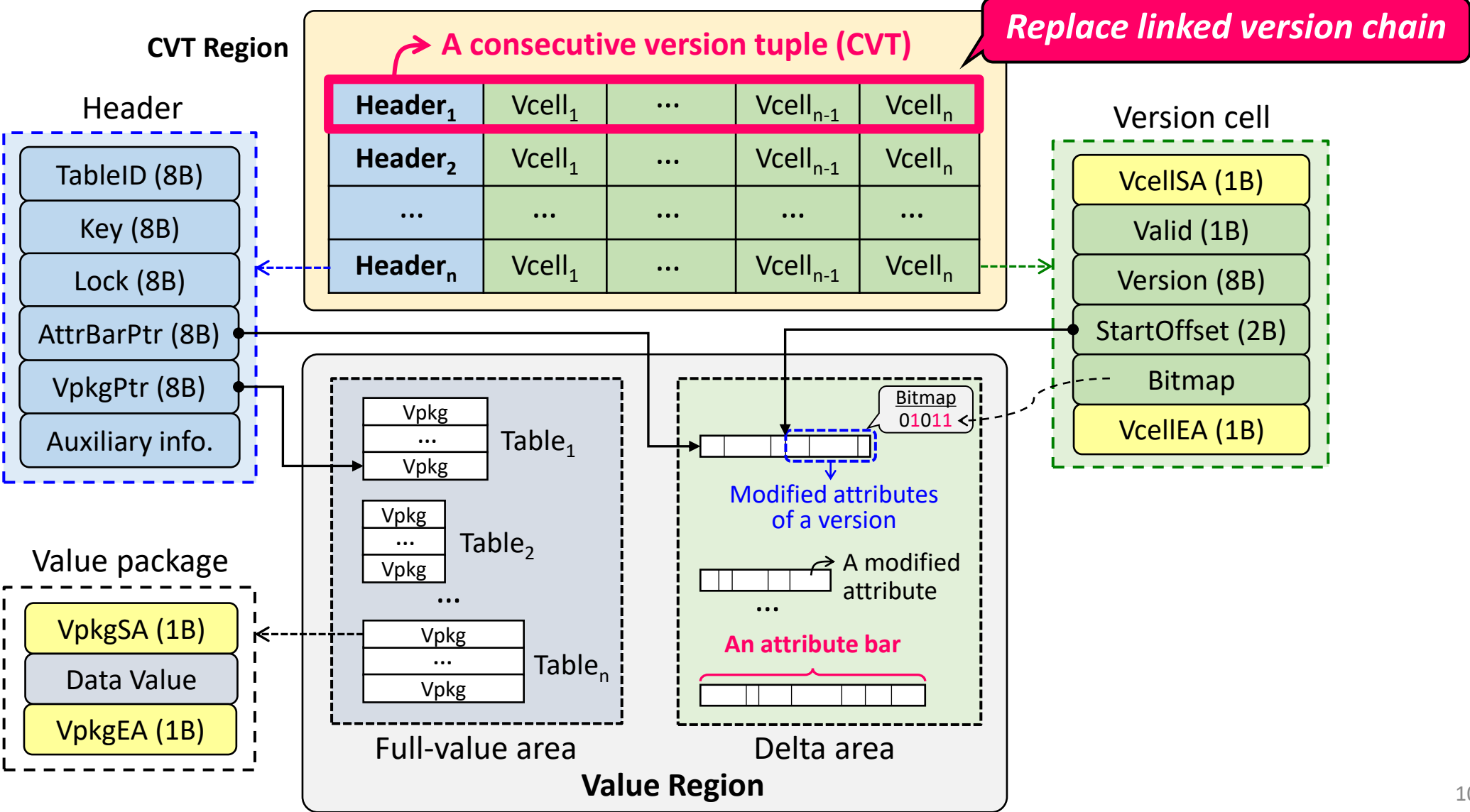




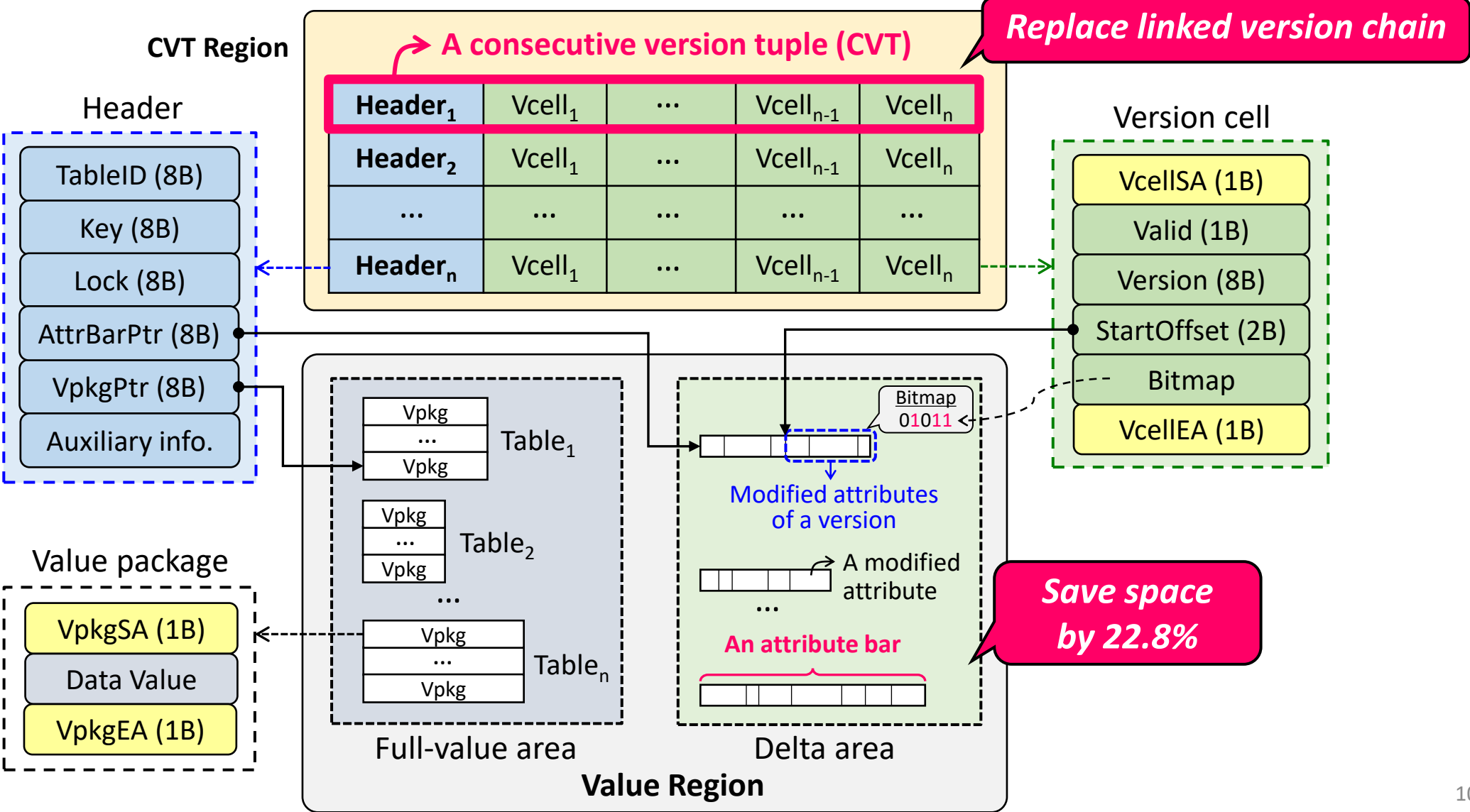
# Consecutive Version Tuple



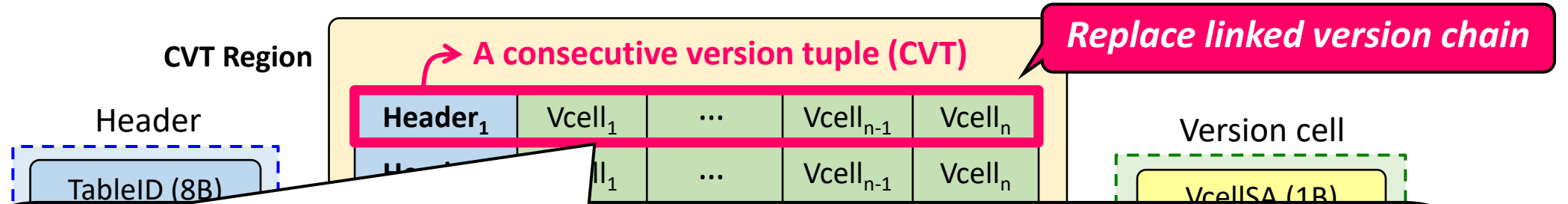
# Consecutive Version Tuple



# Consecutive Version Tuple

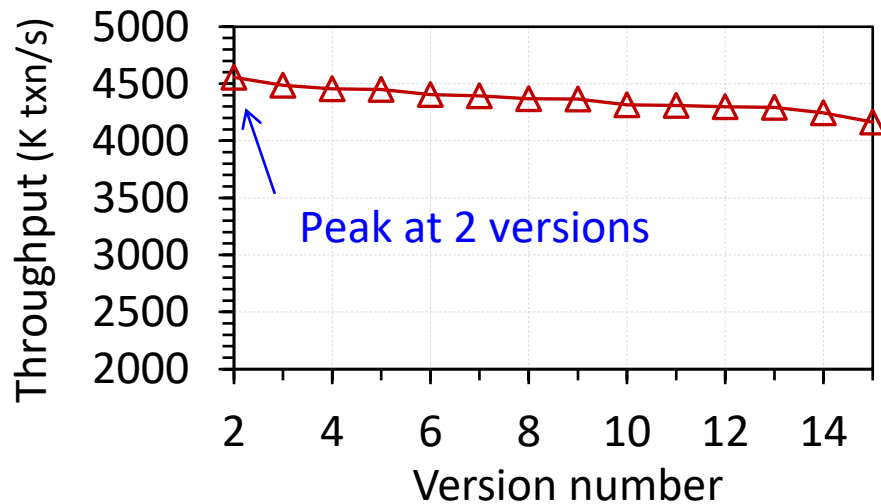


# Consecutive Version Tuple

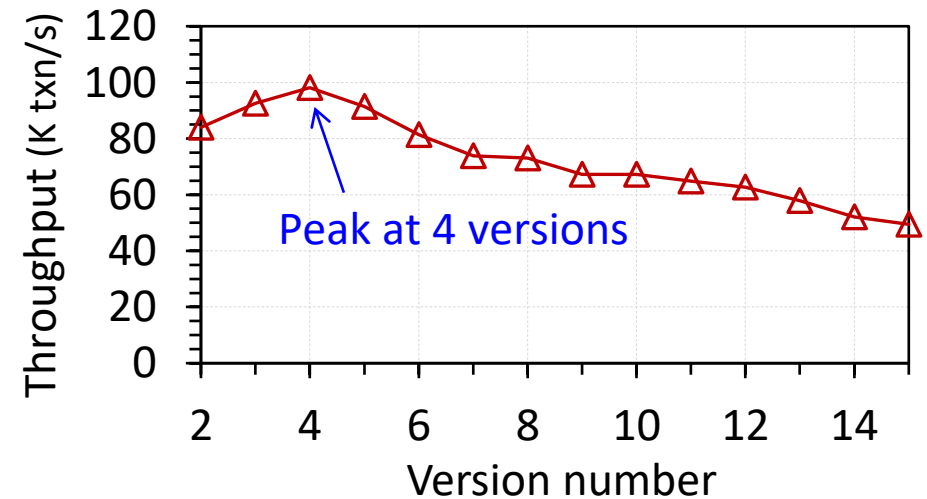


**Number of Versions:** depending on workload characteristics

- Read-write contention
- Number of accessed records



TATP (low contention, short txns)



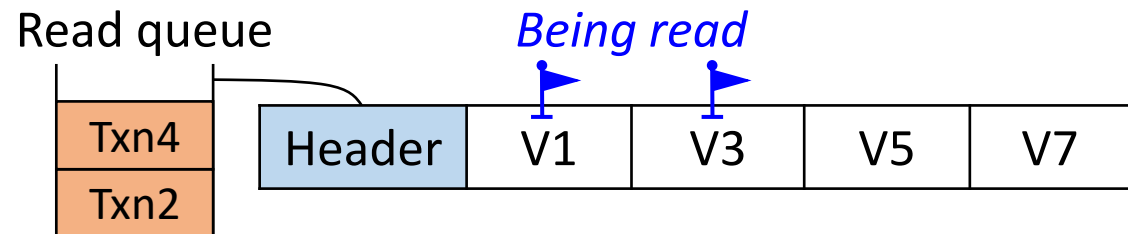
TPCC (high contention, long txns)

# Coordinator-Active Garbage Collection

- A CVT runs out of space – GC required
- Prior systems track transaction states<sup>[1-2]</sup>
  - CPU in memory nodes is too weak to frequently track

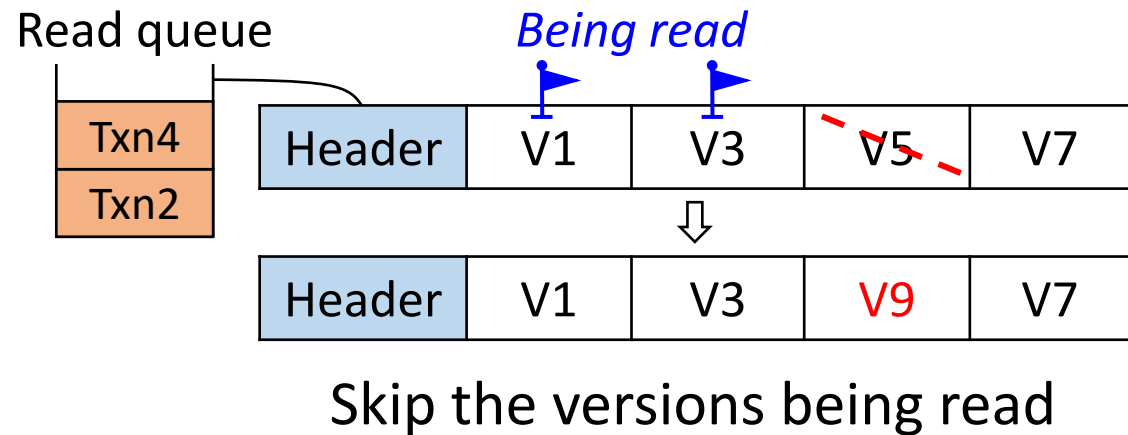
# Coordinator-Active Garbage Collection

- A CVT runs out of space – GC required
- Prior systems track transaction states<sup>[1-2]</sup>
  - CPU in memory nodes is too weak to frequently track



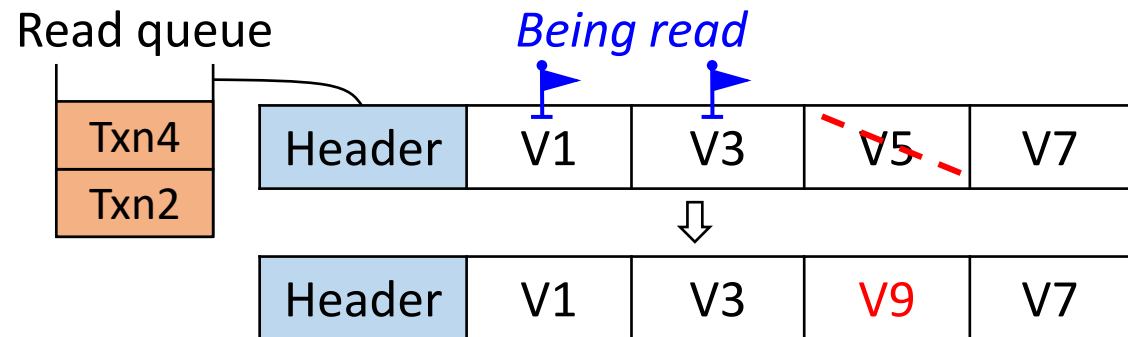
# Coordinator-Active Garbage Collection

- A CVT runs out of space – GC required
- Prior systems track transaction states<sup>[1-2]</sup>
  - CPU in memory nodes is too weak to frequently track



# Coordinator-Active Garbage Collection

- A CVT runs out of space – GC required
- Prior systems track transaction states<sup>[1-2]</sup>
  - CPU in memory nodes is too weak to frequently track



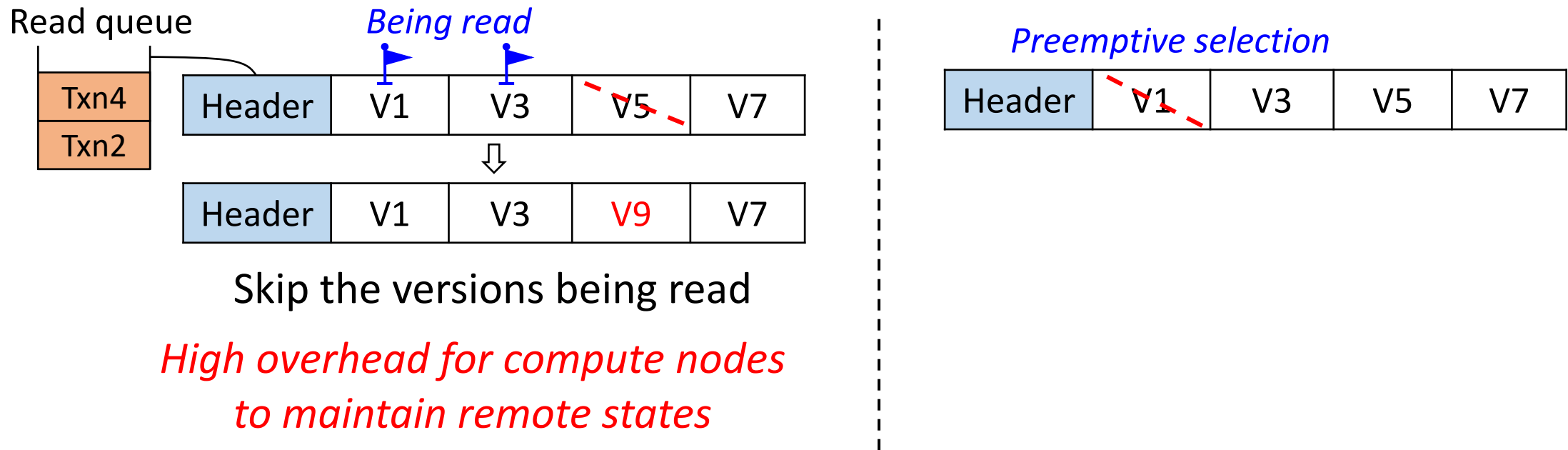
Skip the versions being read

*High overhead for compute nodes  
to maintain remote states*



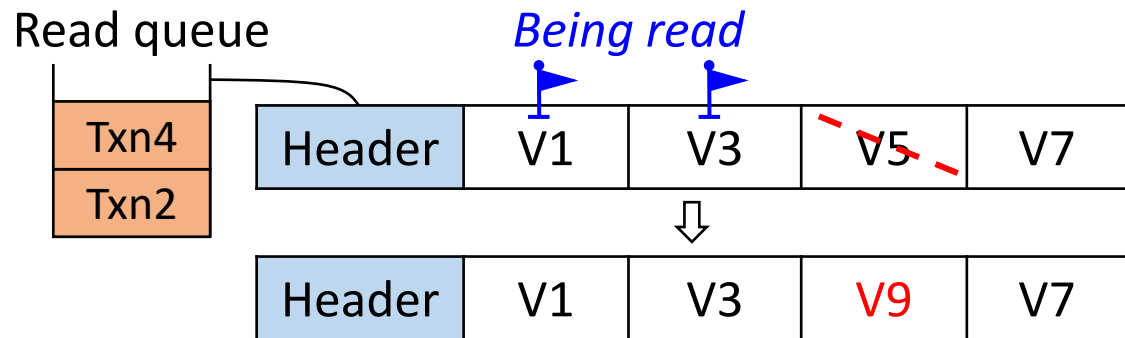
# Coordinator-Active Garbage Collection

- A CVT runs out of space – GC required
- Prior systems track transaction states<sup>[1-2]</sup>
  - CPU in memory nodes is too weak to frequently track



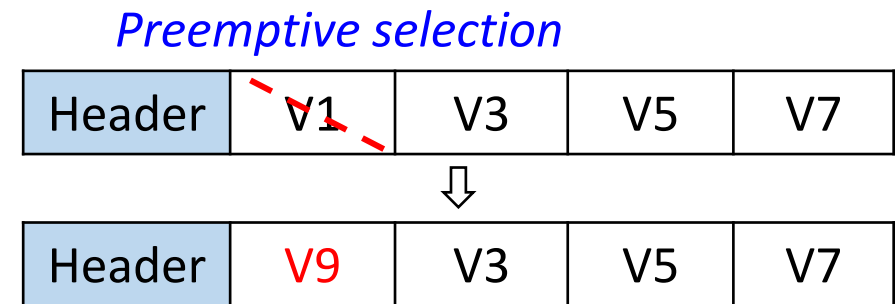
# Coordinator-Active Garbage Collection

- A CVT runs out of space – GC required
- Prior systems track transaction states<sup>[1-2]</sup>
  - CPU in memory nodes is too weak to frequently track



Skip the versions being read

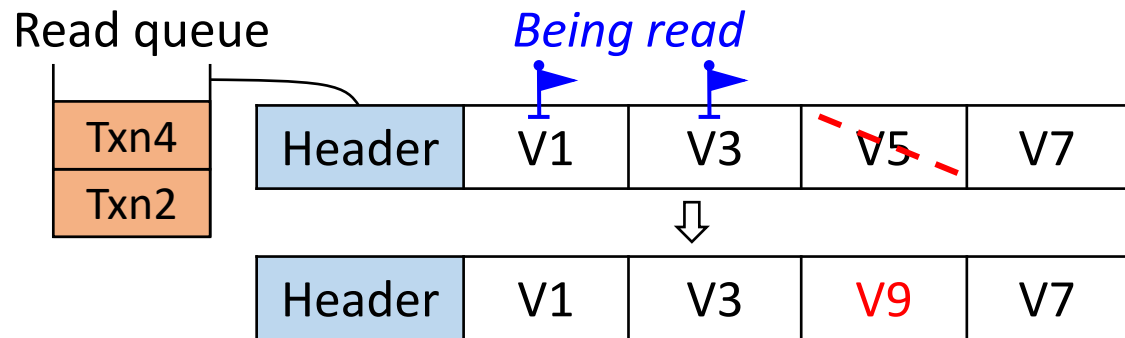
*High overhead for compute nodes  
to maintain remote states*



Overwrite the oldest version

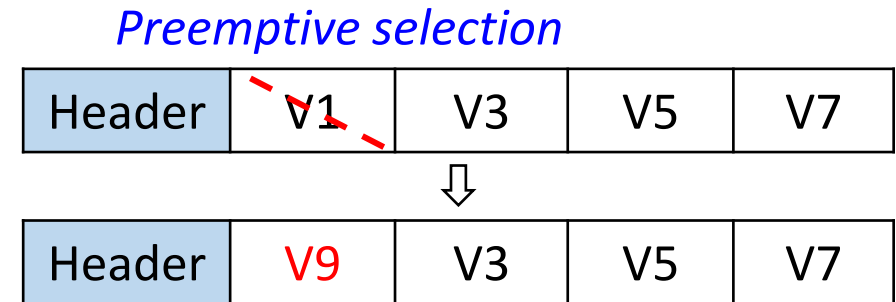
# Coordinator-Active Garbage Collection

- A CVT runs out of space – GC required
- Prior systems track transaction states<sup>[1-2]</sup>
  - CPU in memory nodes is too weak to frequently track



Skip the versions being read

*High overhead for compute nodes  
to maintain remote states*

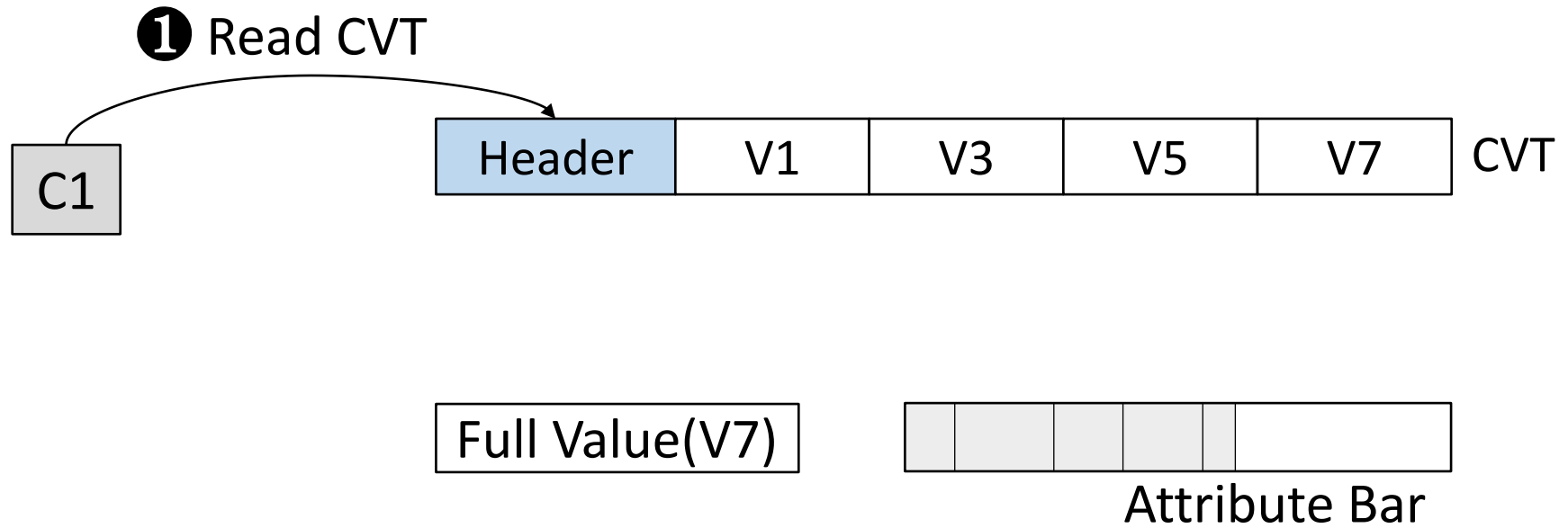


Overwrite the oldest version

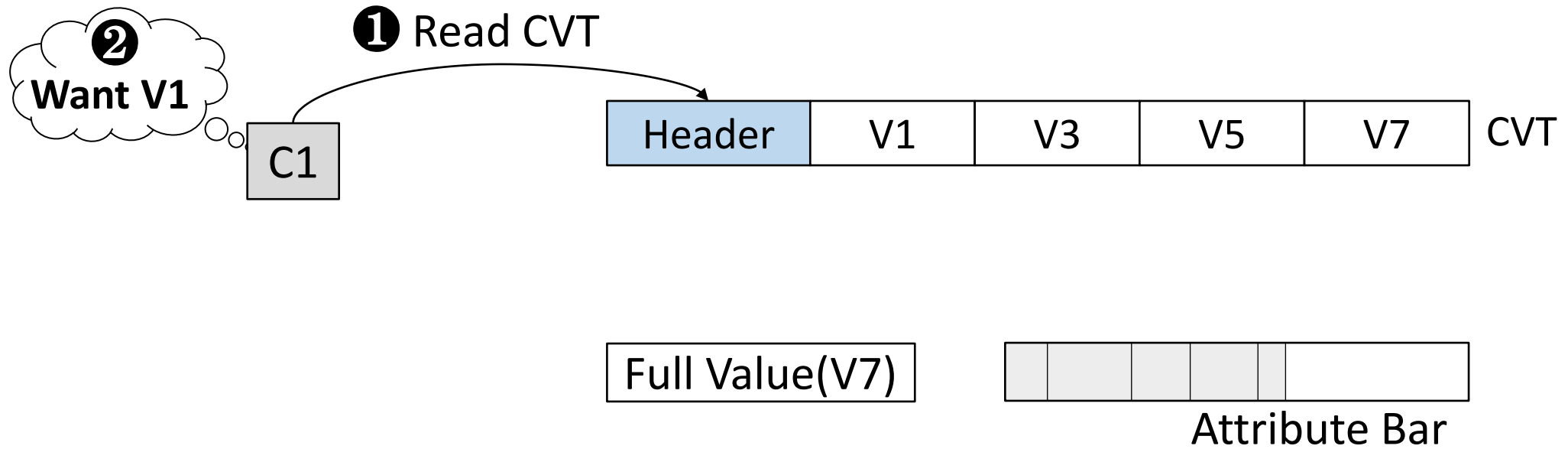
*Simple, no tracking  
Low abort rate with fast RDMA*

# Anchor-Assisted Read

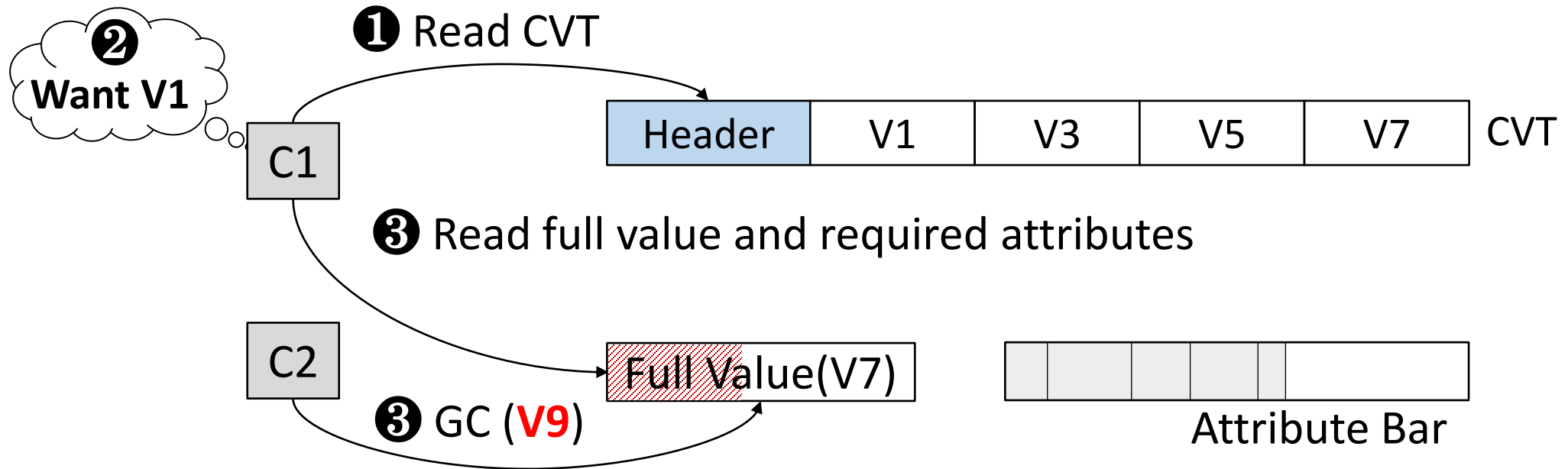
# Anchor-Assisted Read



# Anchor-Assisted Read



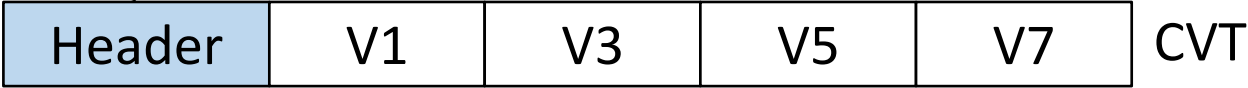
# Anchor-Assisted Read



# Anchor-Assisted Read

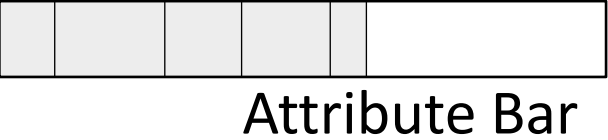
②  
Want V1

① Read CVT



③ Read full value and required attributes  
(may read a corrupted value)

C2

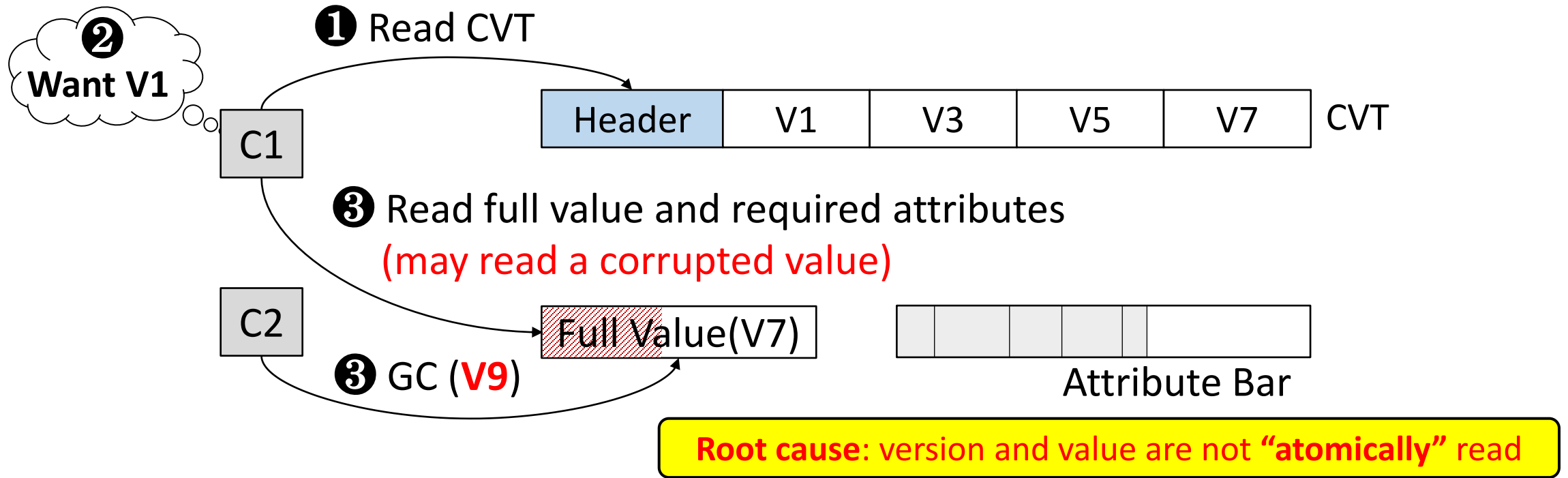


③ GC (V9)

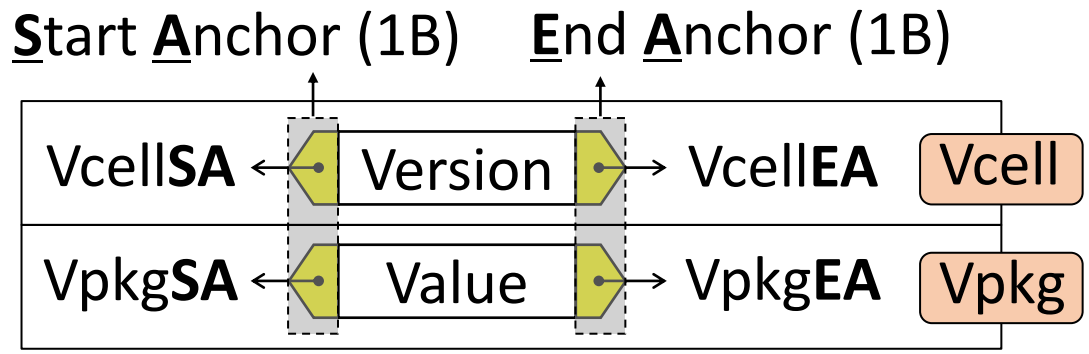
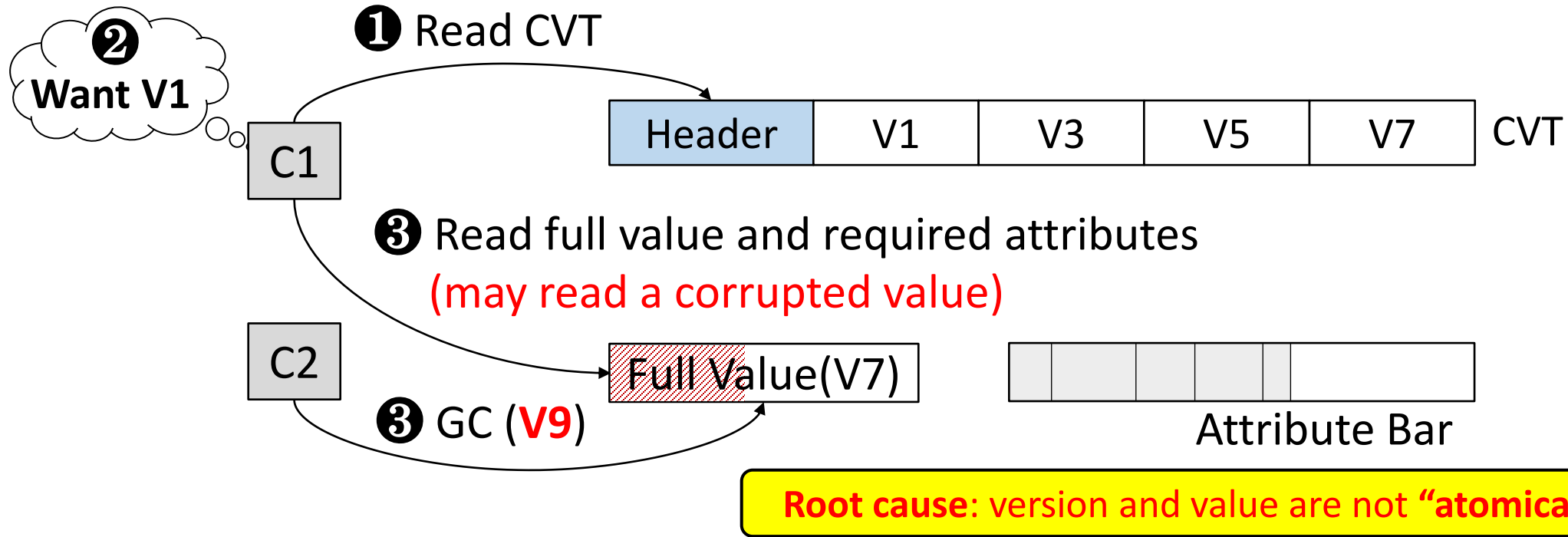
C1



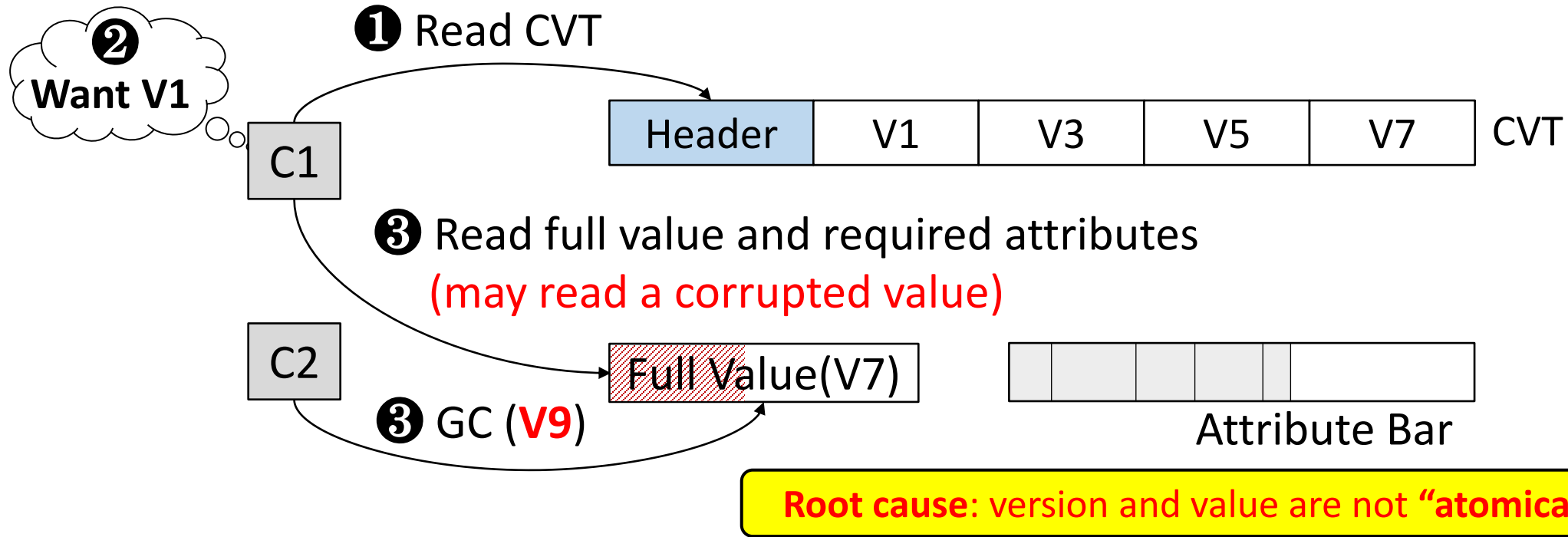
# Anchor-Assisted Read



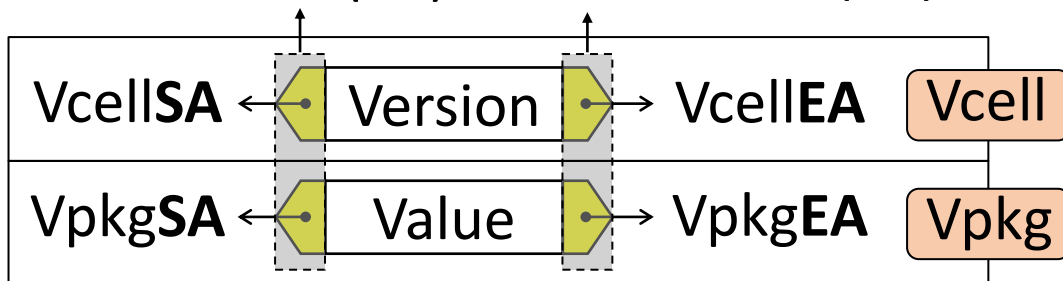
# Anchor-Assisted Read



# Anchor-Assisted Read

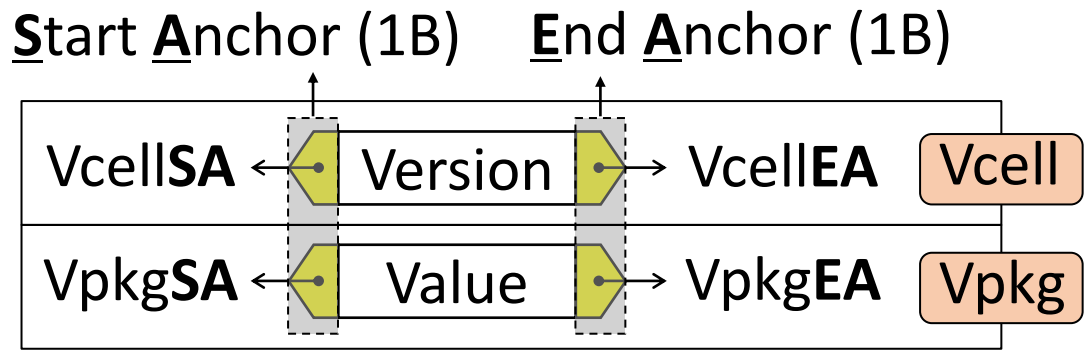
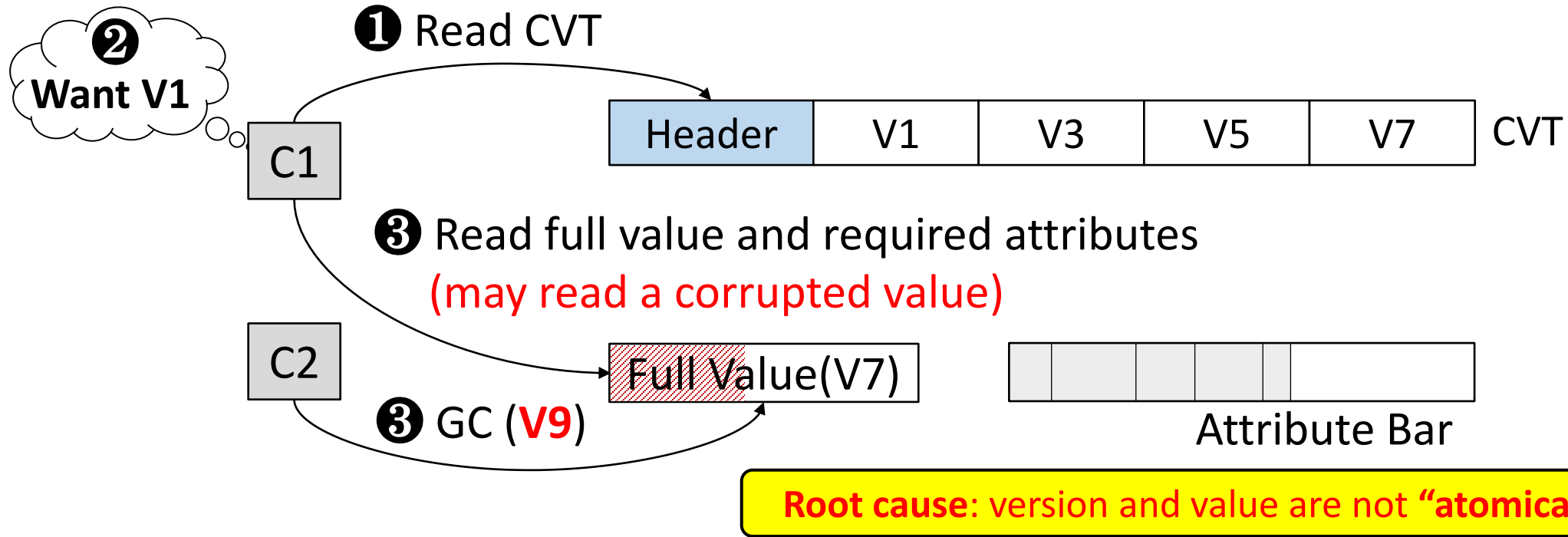


Start Anchor (1B)    End Anchor (1B)



• **Writer:** Vpkg → attributes → Vcell

# Anchor-Assisted Read



- **Writer:** Vpkg → attributes → Vcell
  - **Reader:** check "all anchors are equal"
- VcellSA = VcellEA = VpkgSA = VpkgEA** ✓

# One-Sided RDMA-Based MVCC

**A** CVT   **A** Vpkg and any required attributes   **A** Batched writes    Lock    Unlock    RTT

*Coordinator*

A's primary

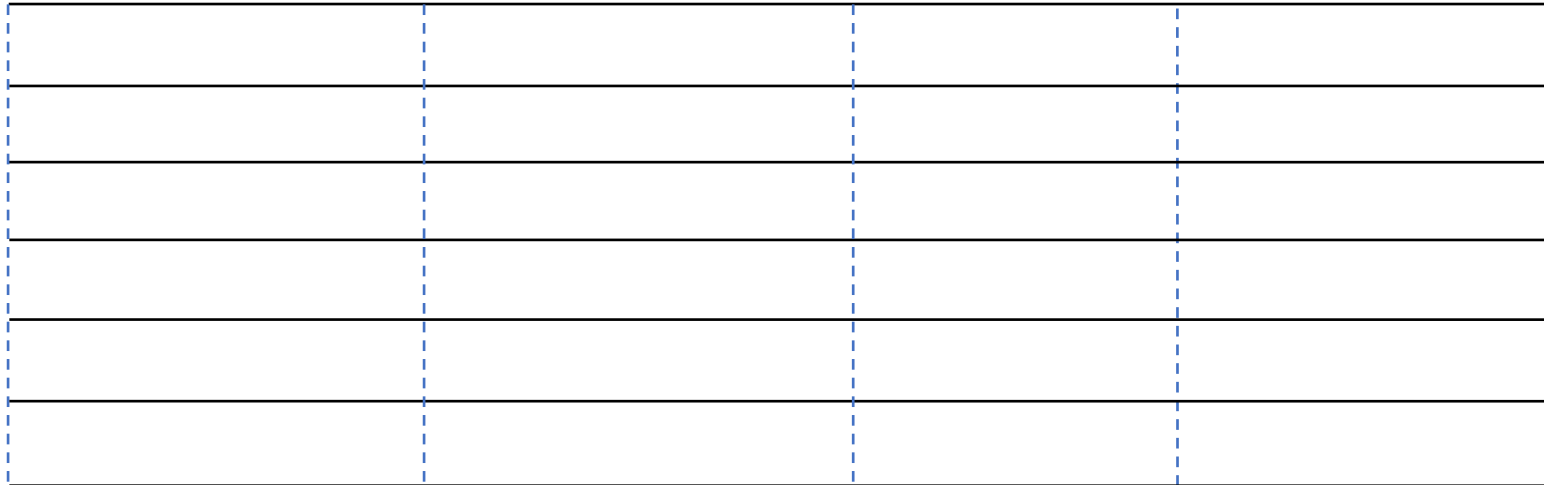
A's backups

B's primary

B's backups

C's primary

C's backups



***Txn begin***

Read A,B,C

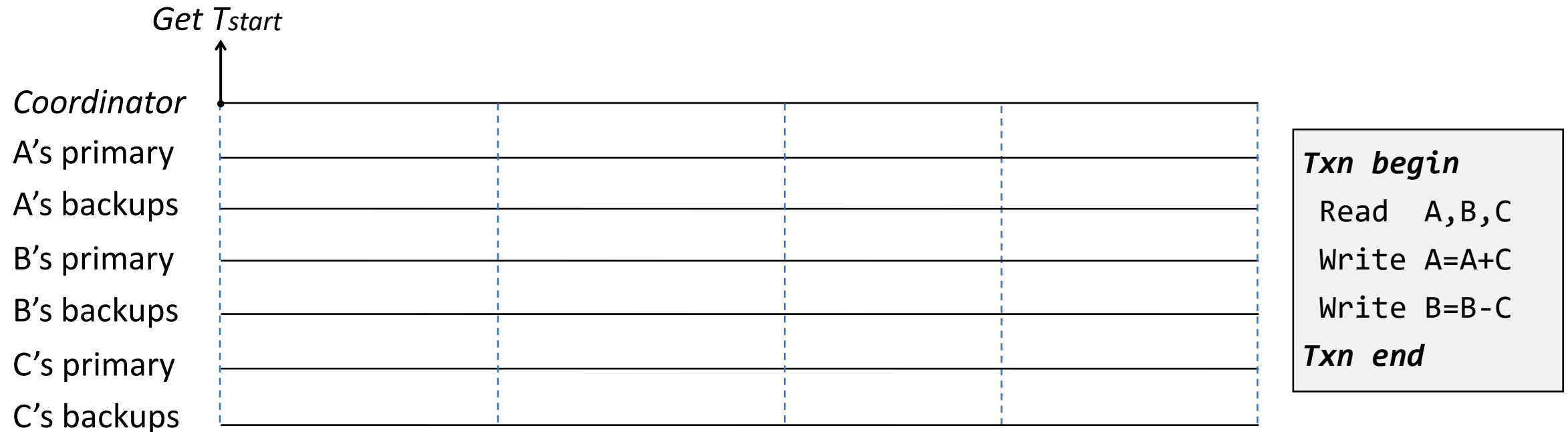
Write  $A=A+C$

Write  $B=B-C$

***Txn end***

# One-Sided RDMA-Based MVCC

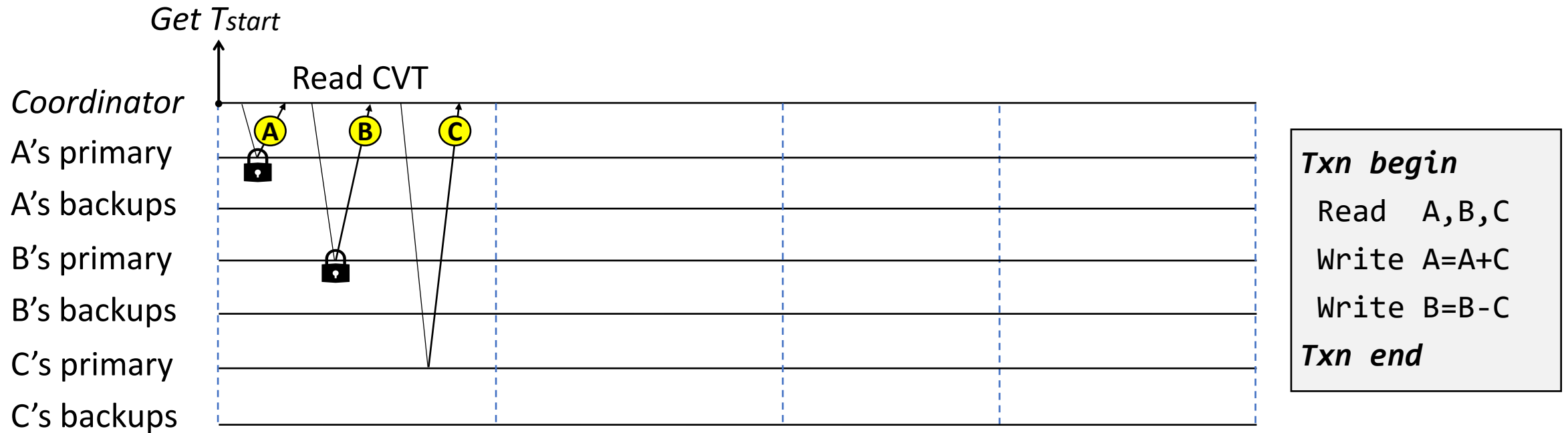
**A** CVT   **A** Vpkg and any required attributes   **A** Batched writes    Lock    Unlock    RTT



\* Tstart/Tcommit stands for start/commit timestamp

# One-Sided RDMA-Based MVCC

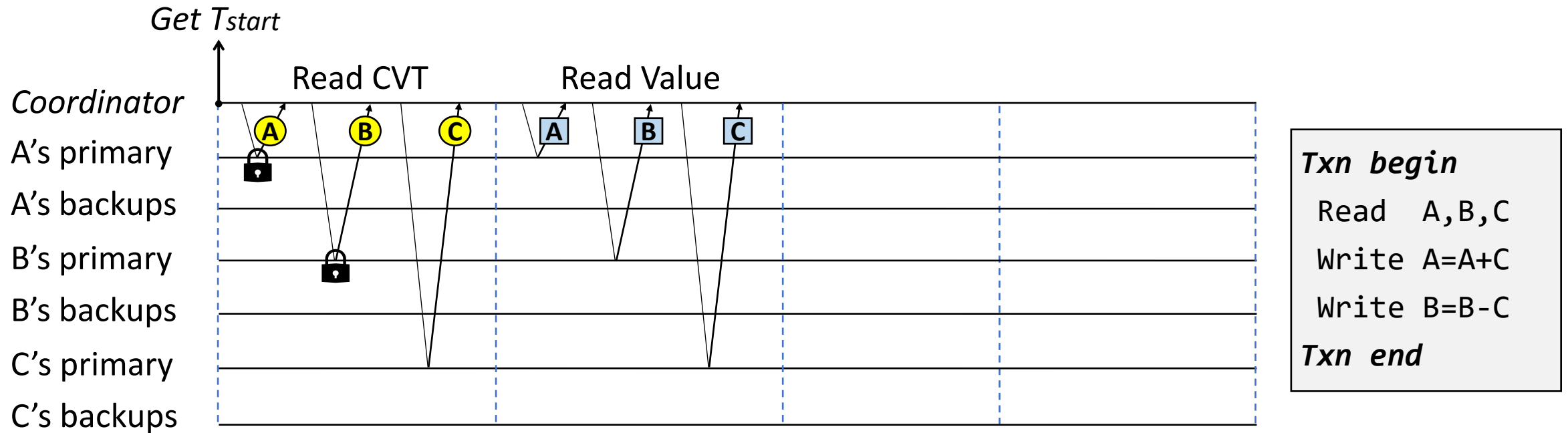
A CVT   
 A Vpkg and any required attributes   
 A Batched writes   
  Lock   
  Unlock   
  RTT



\* T<sub>start</sub>/T<sub>commit</sub> stands for start/commit timestamp

# One-Sided RDMA-Based MVCC

A CVT   
 A Vpkg and any required attributes   
 A Batched writes   
 🔒 Lock   
 🔓 Unlock   
 ↔ RTT

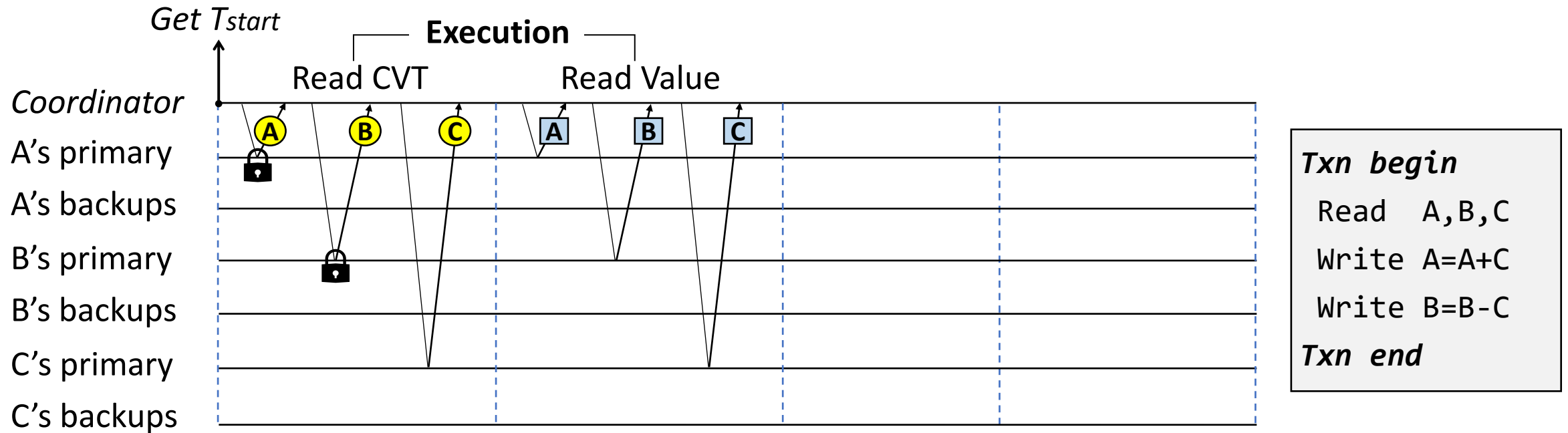


\* Tstart/Tcommit stands for start/commit timestamp



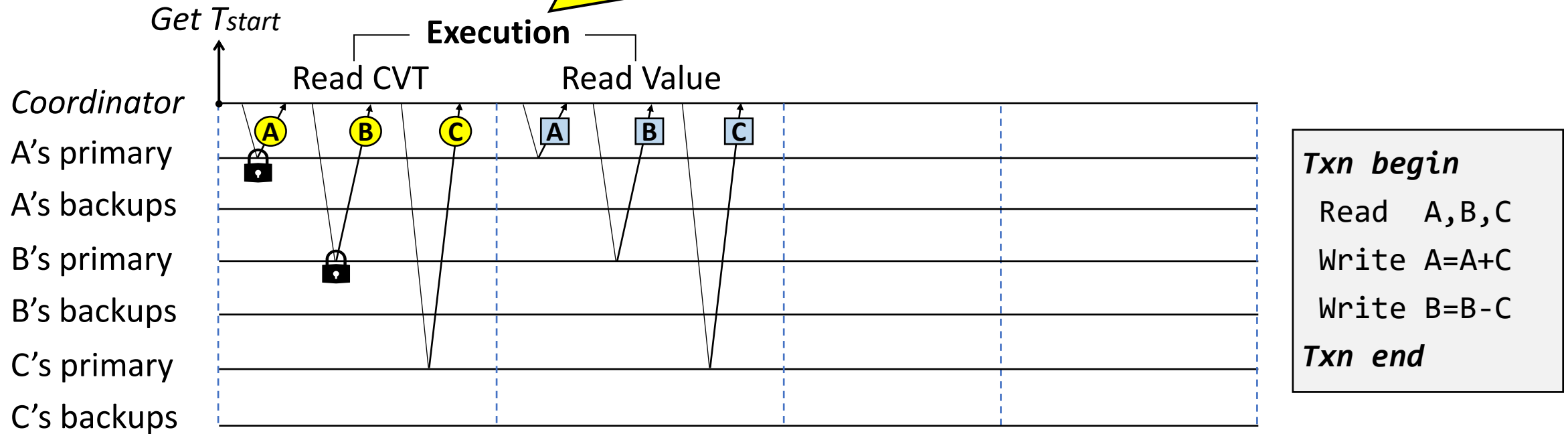
# One-Sided RDMA-Based MVCC

A CVT   
 A Vpkg and any required attributes   
 A Batched writes   
 🔒 Lock   
 🔓 Unlock   
 ↔ RTT



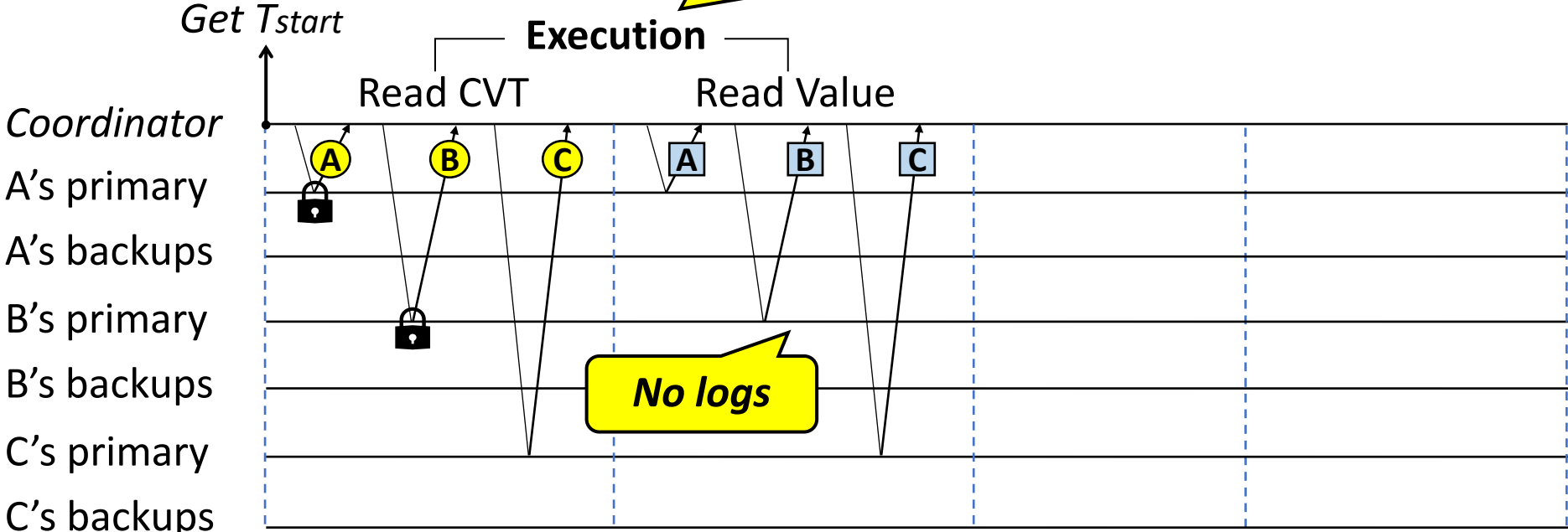
\* T<sub>start</sub>/T<sub>commit</sub> stands for start/commit timestamp

# One-Sided RDMA-Based MVCC



\*  $T_{start}/T_{commit}$  stands for start/commit timestamp

# One-Sided RDMA-Based MVCC



*Txn begin*

Read A,B,C

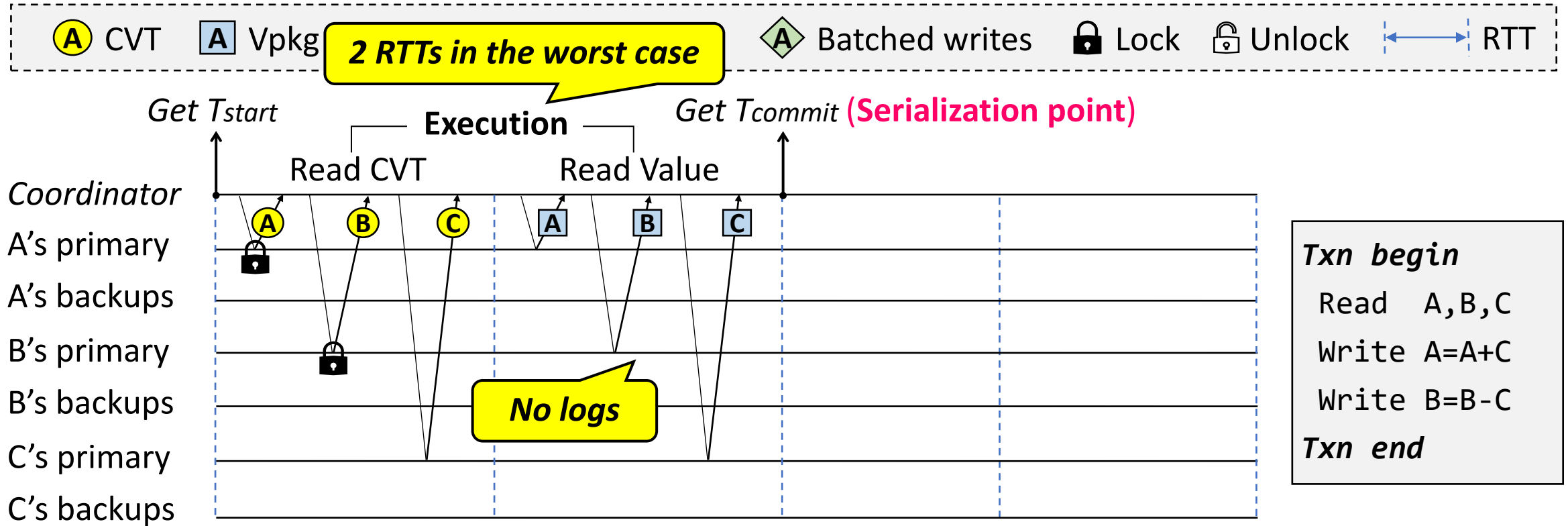
Write  $A=A+C$

Write  $B=B-C$

*Txn end*

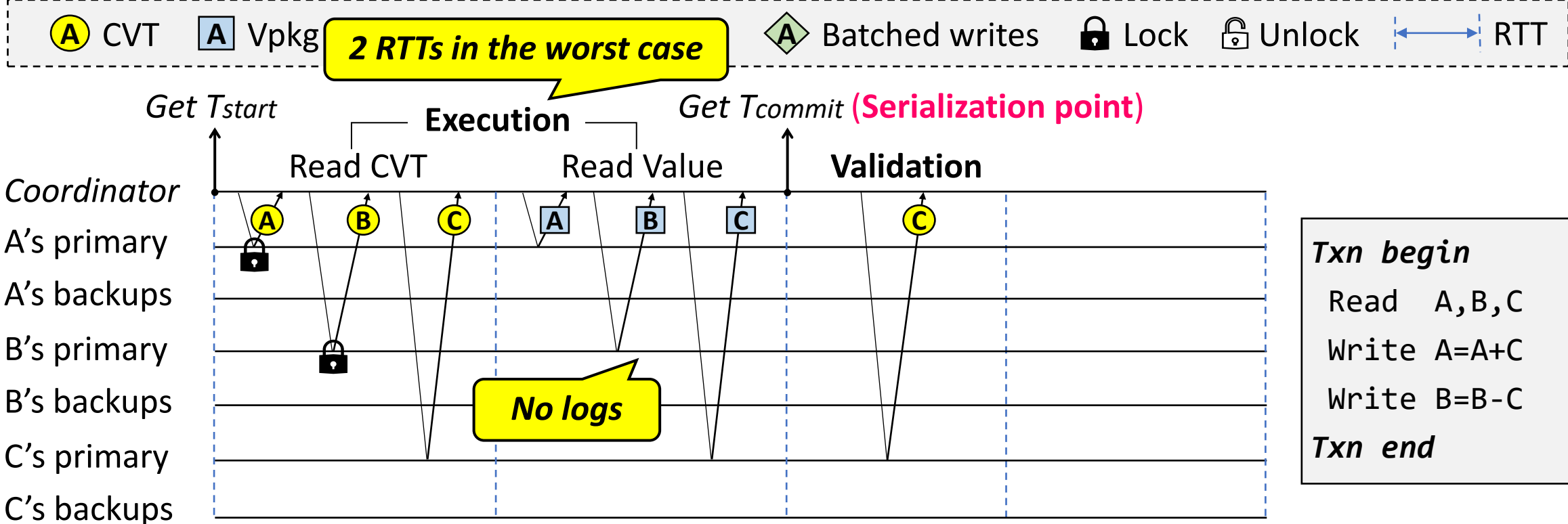
\*  $T_{start}/T_{commit}$  stands for start/commit timestamp

# One-Sided RDMA-Based MVCC



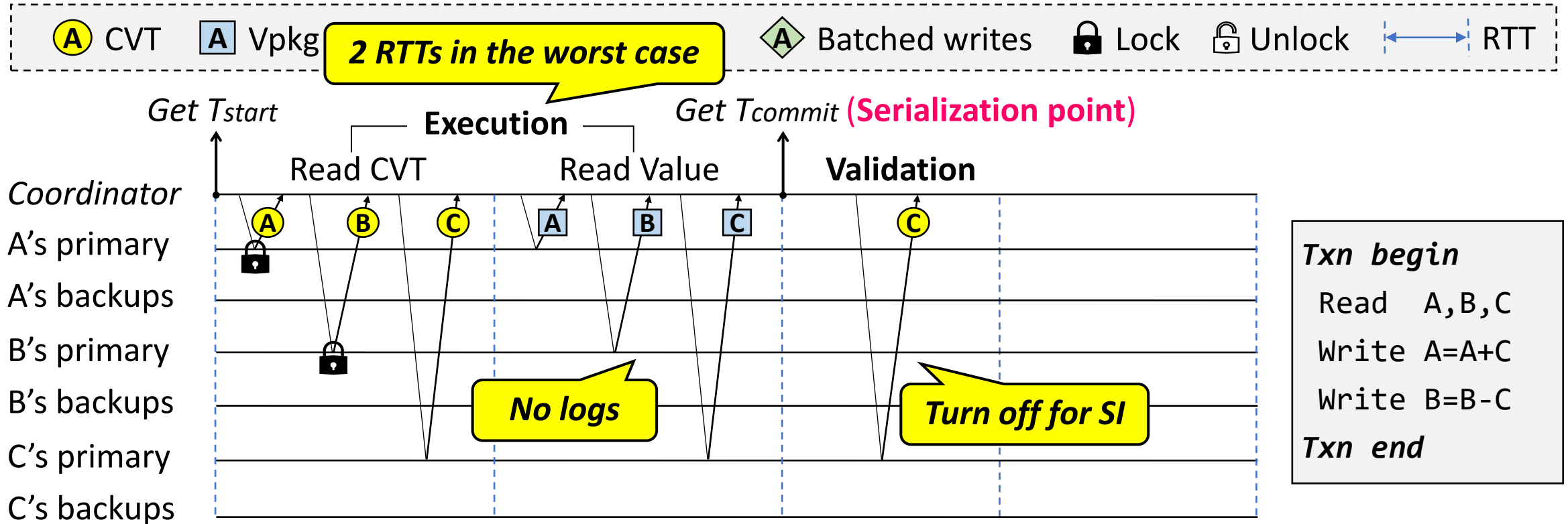
\*  $T_{start}/T_{commit}$  stands for start/commit timestamp

# One-Sided RDMA-Based MVCC



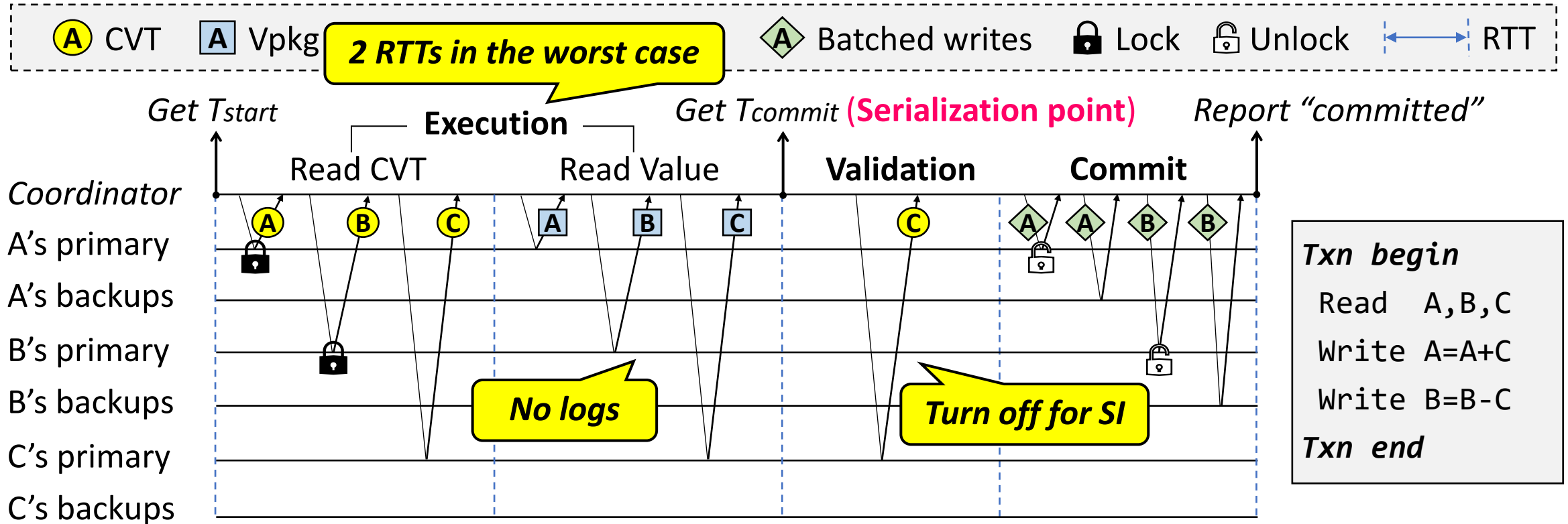
\* Tstart/Tcommit stands for start/commit timestamp

# One-Sided RDMA-Based MVCC



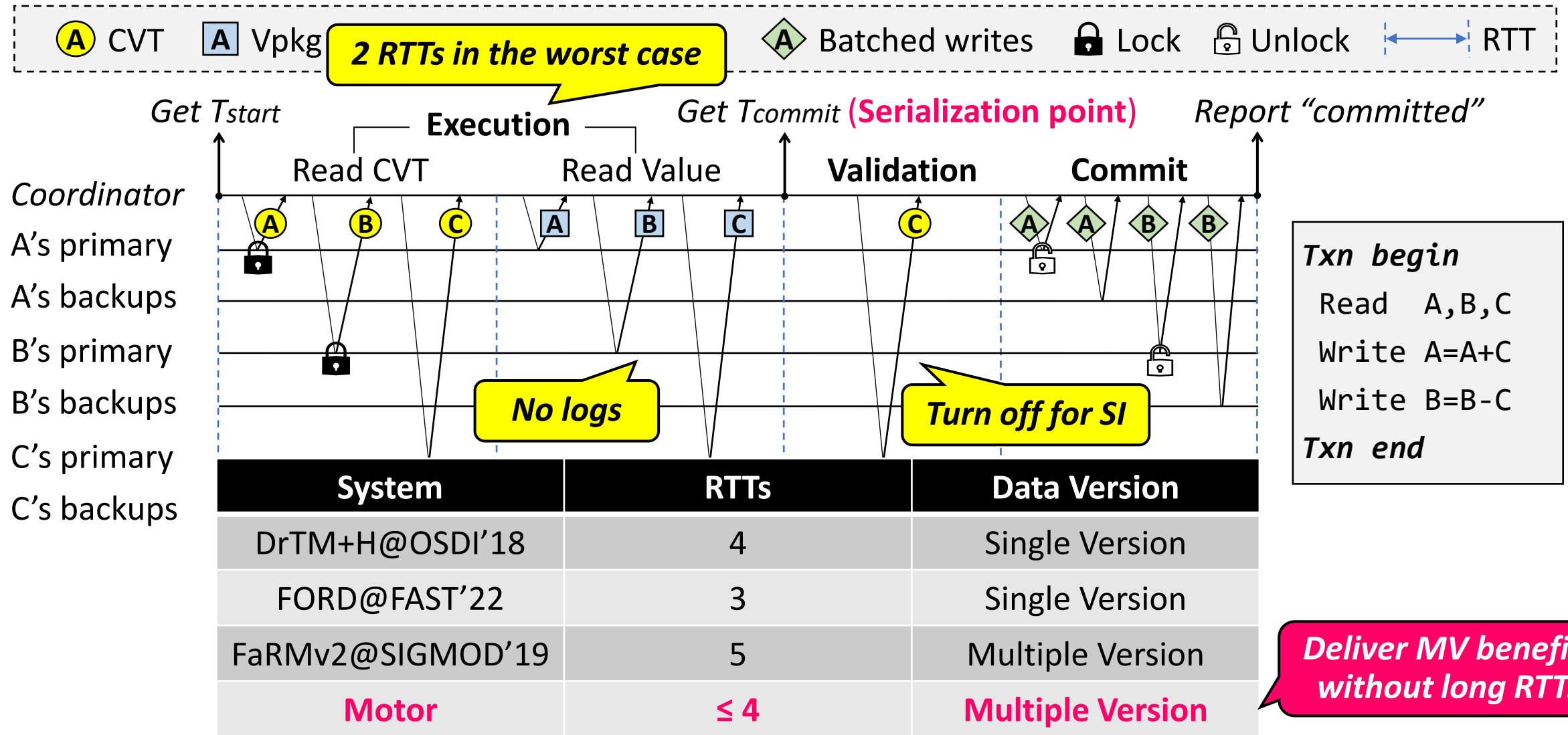
\* T<sub>start</sub>/T<sub>commit</sub> stands for start/commit timestamp

# One-Sided RDMA-Based MVCC



\*  $T_{start}/T_{commit}$  stands for start/commit timestamp

# One-Sided RDMA-Based MVCC



\* Tstart/Tcommit stands for start/commit timestamp



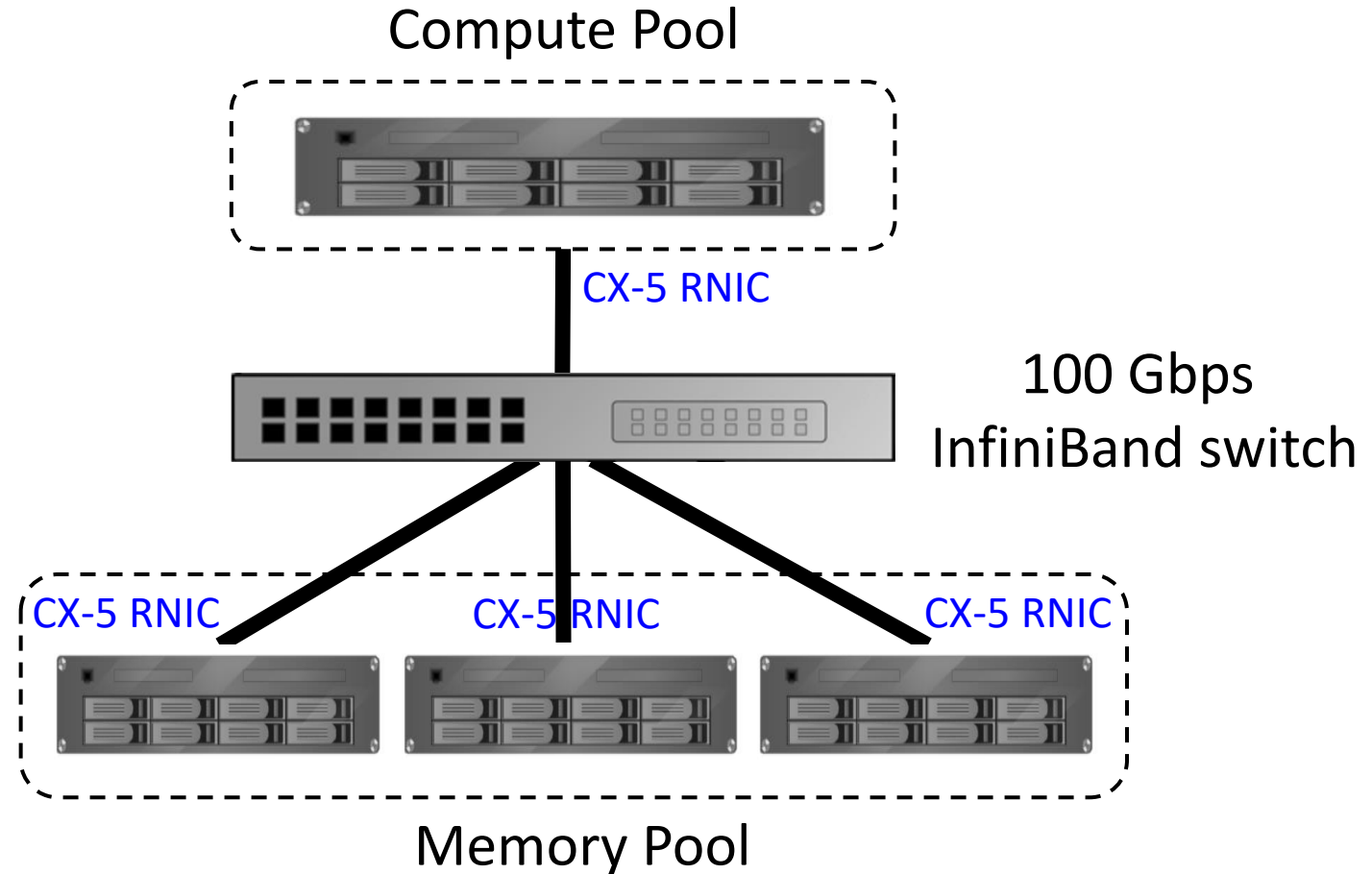
# Evaluation

## ➤ Workloads

- 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

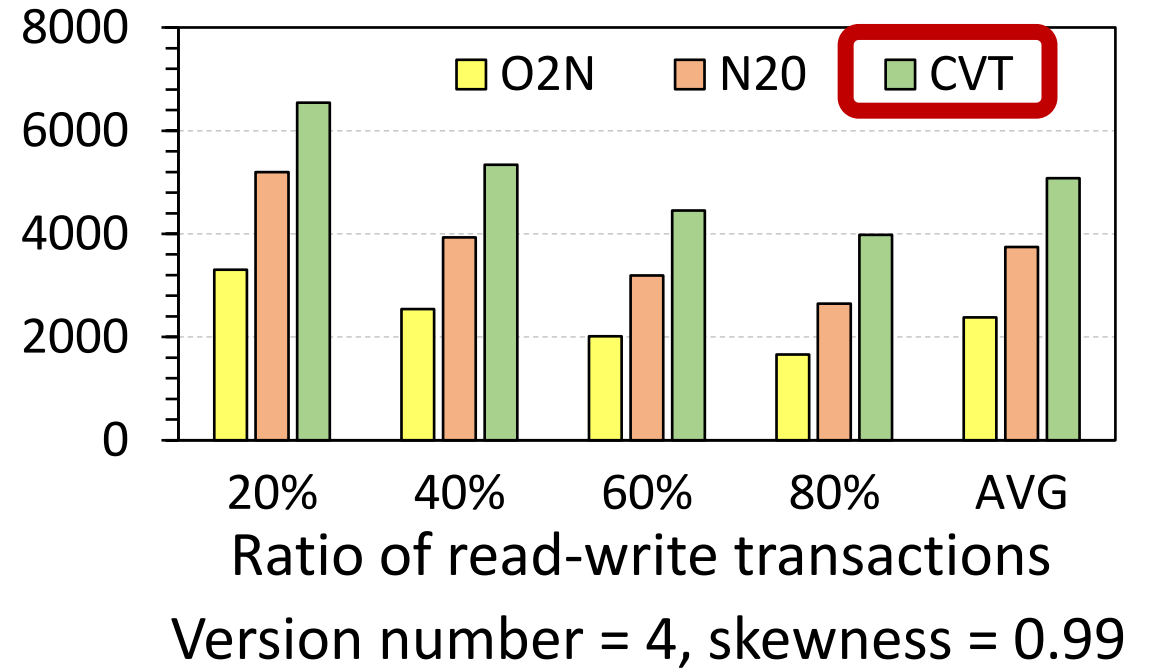
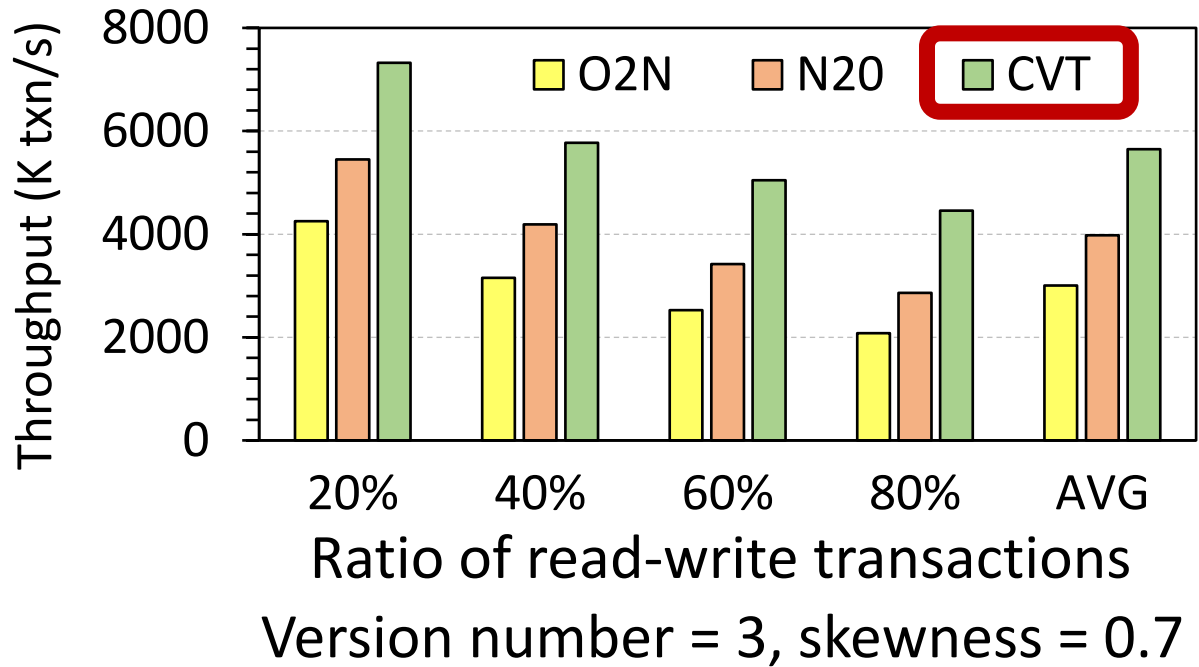
## ➤ Comparisons

- FaRMv2@SIGMOD'19 (referred as FaRMv2-DM)
- FORD@FAST'22



# Performance of Version Structures

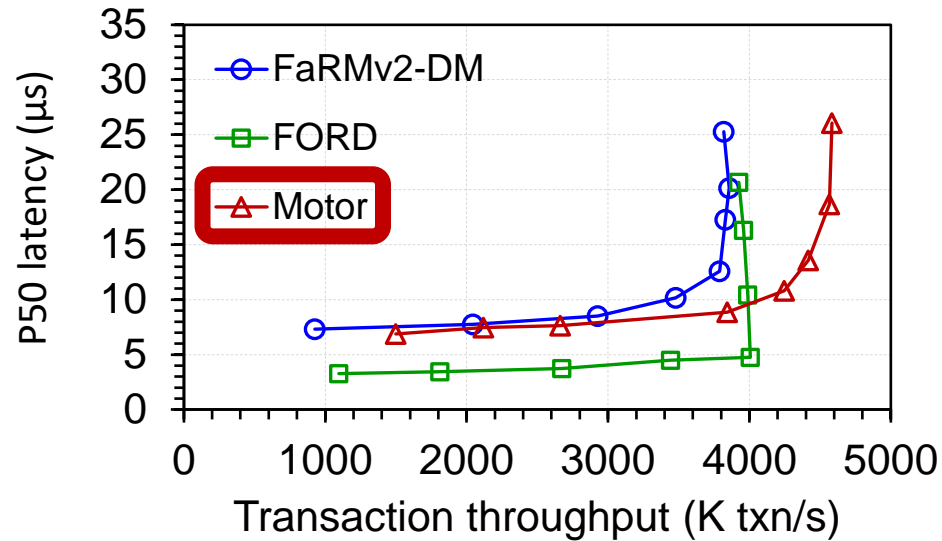
## ➤ KV store



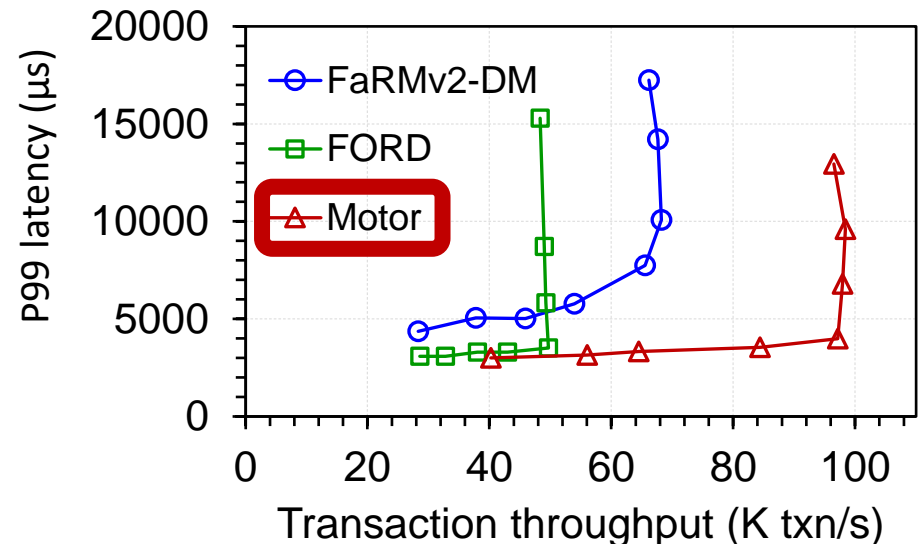
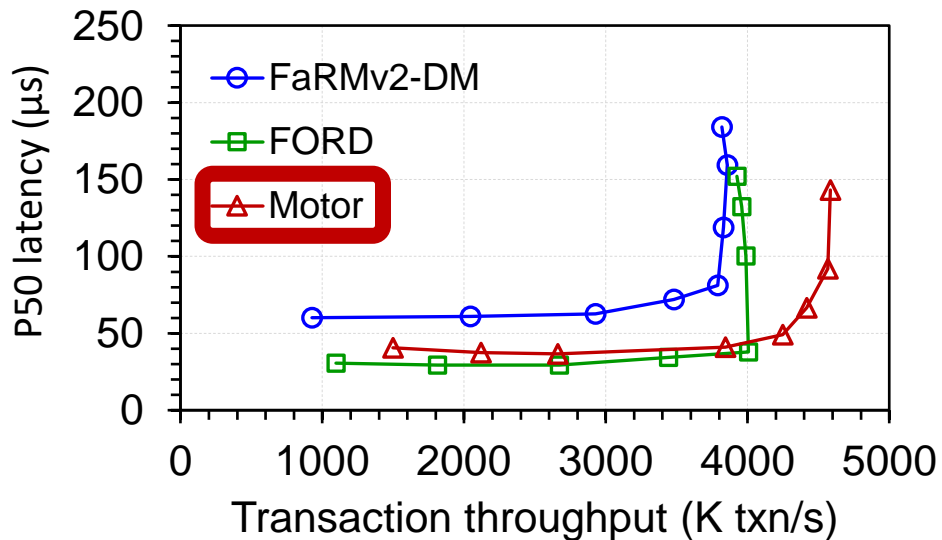
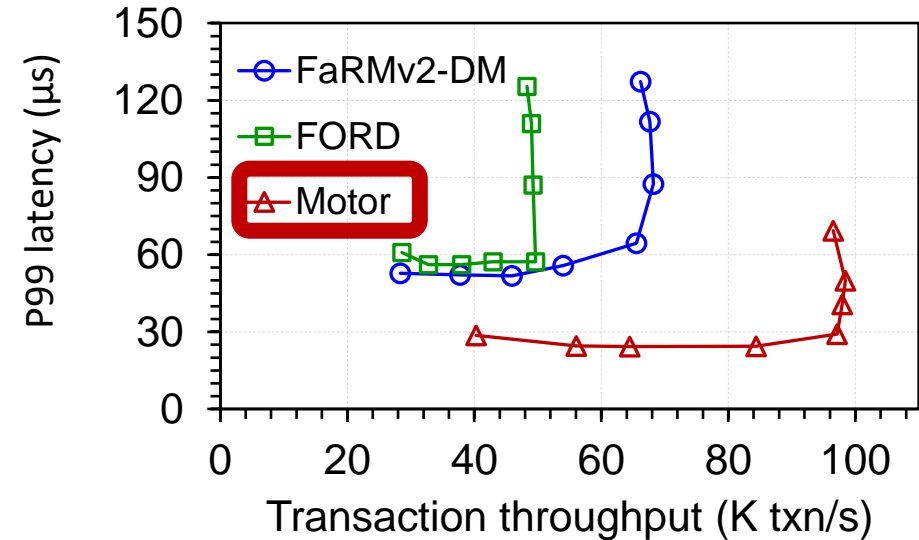
CVT improves throughput by  $\left\langle \begin{array}{l} 1.7-2.4x \text{ over O2N} \\ 1.3-1.6x \text{ over N20} \end{array} \right.$

# End-to-End Performance

TATP (read-intensive)

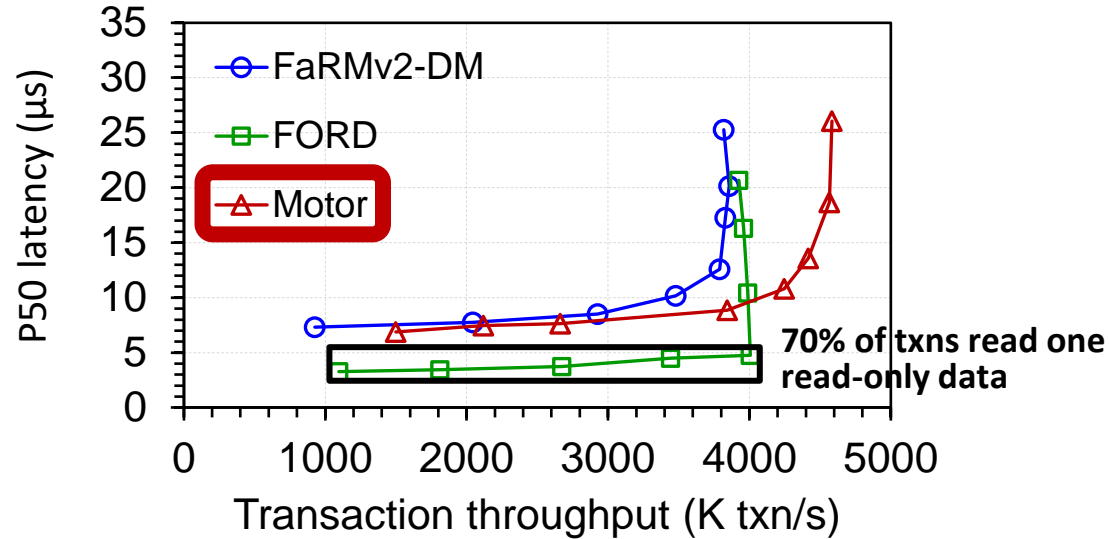


TPCC (write-intensive)

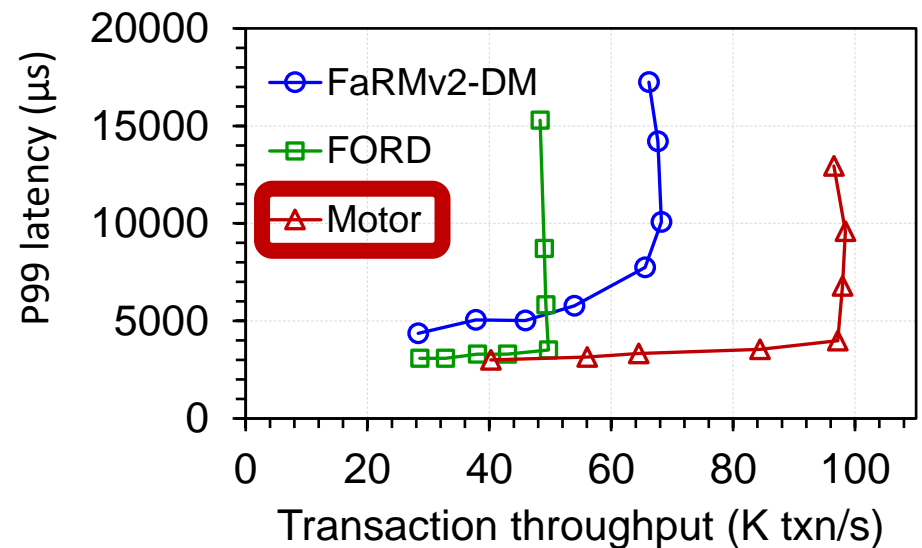
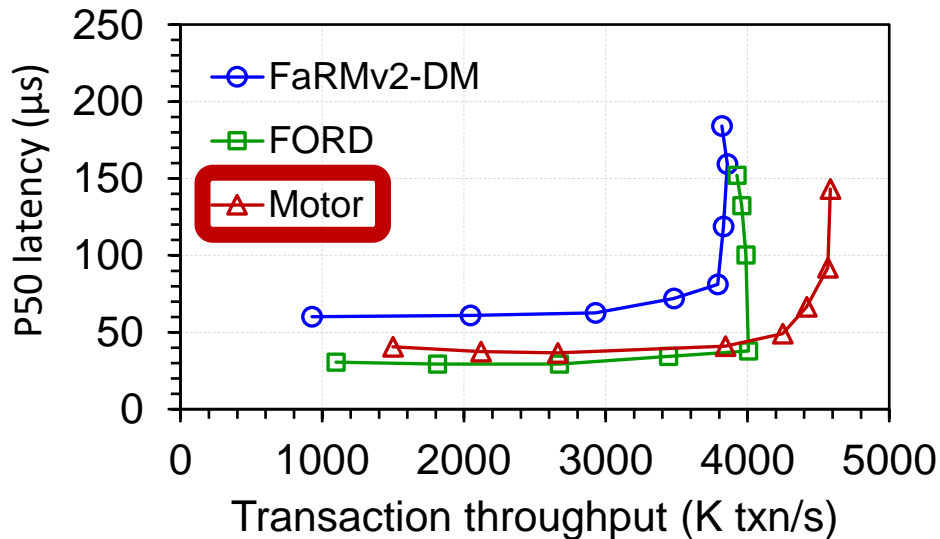
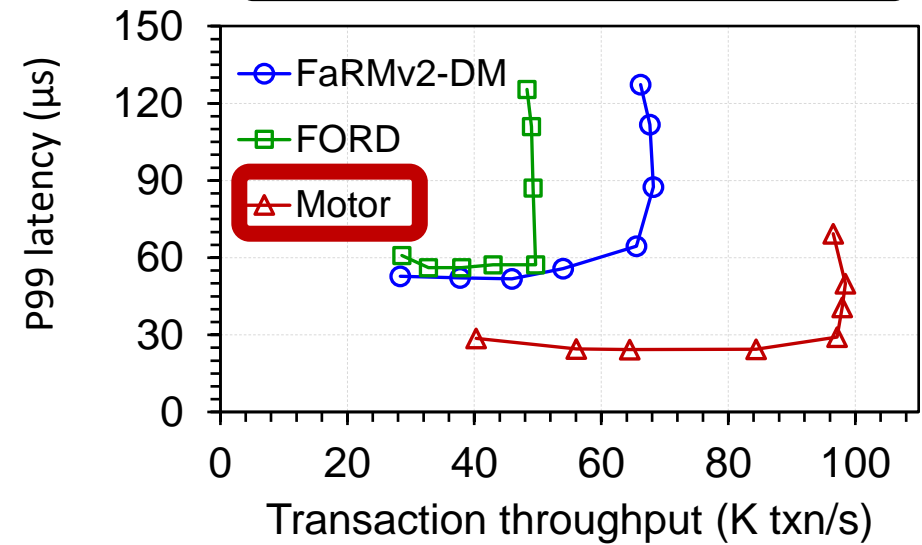


# End-to-End Performance

TATP (read-intensive)

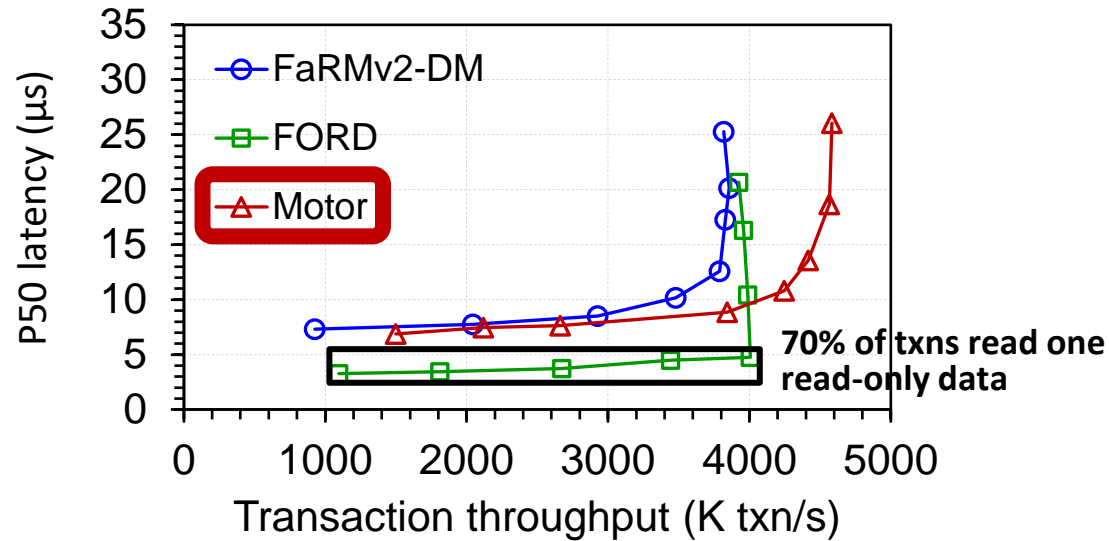


TPCC (write-intensive)

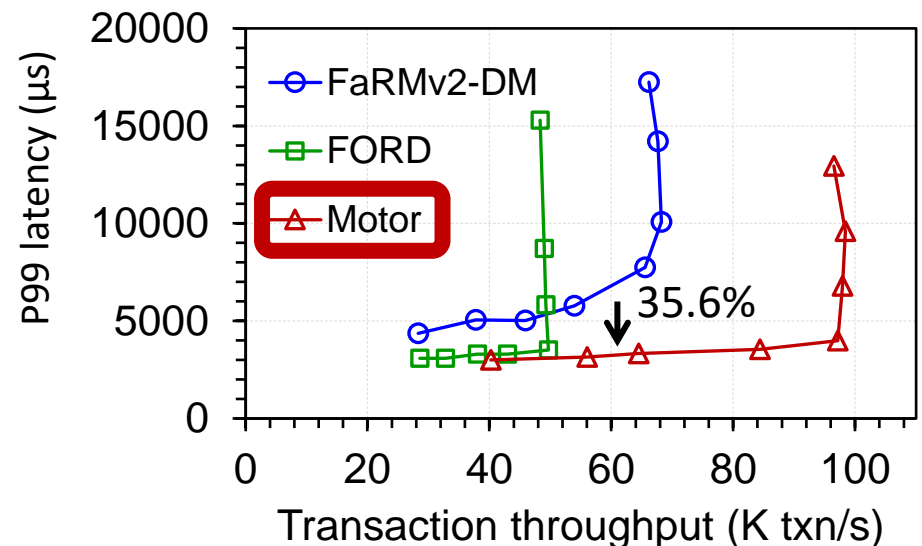
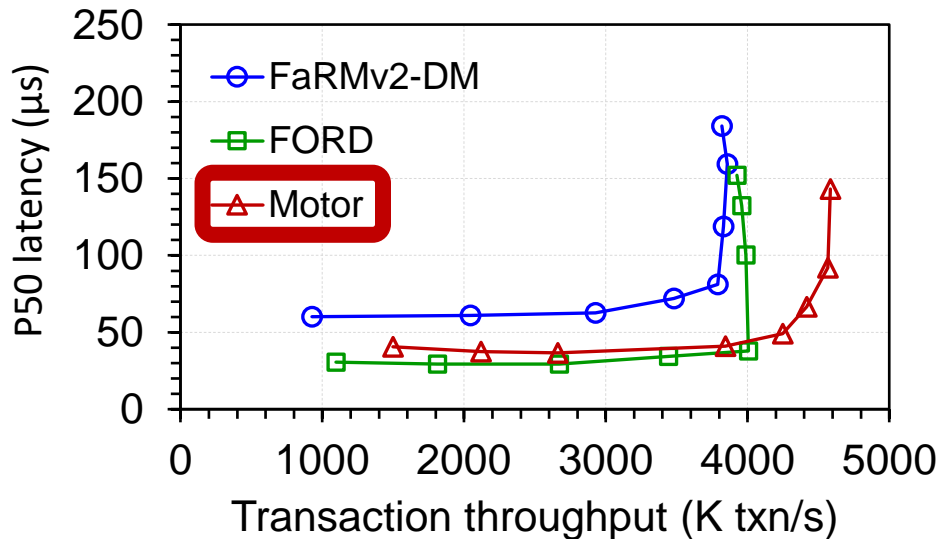
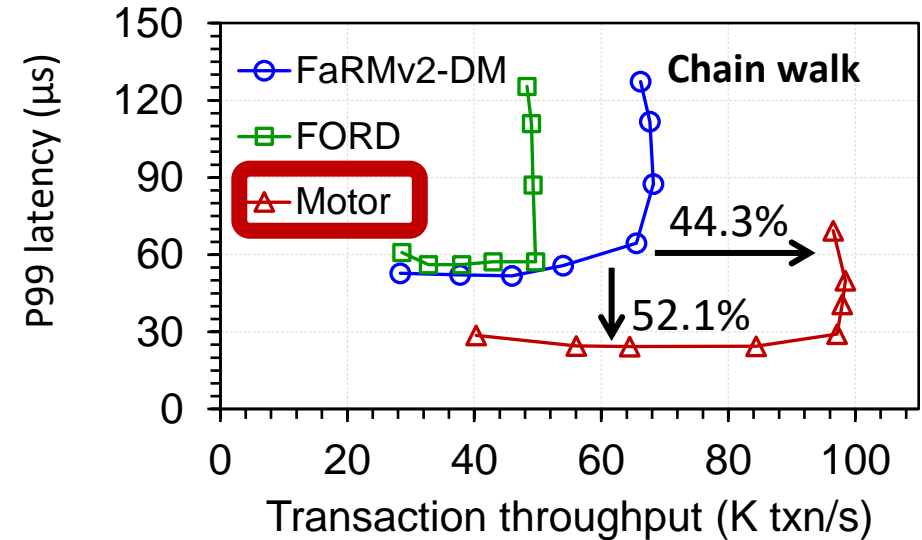


# End-to-End Performance

TATP (read-intensive)

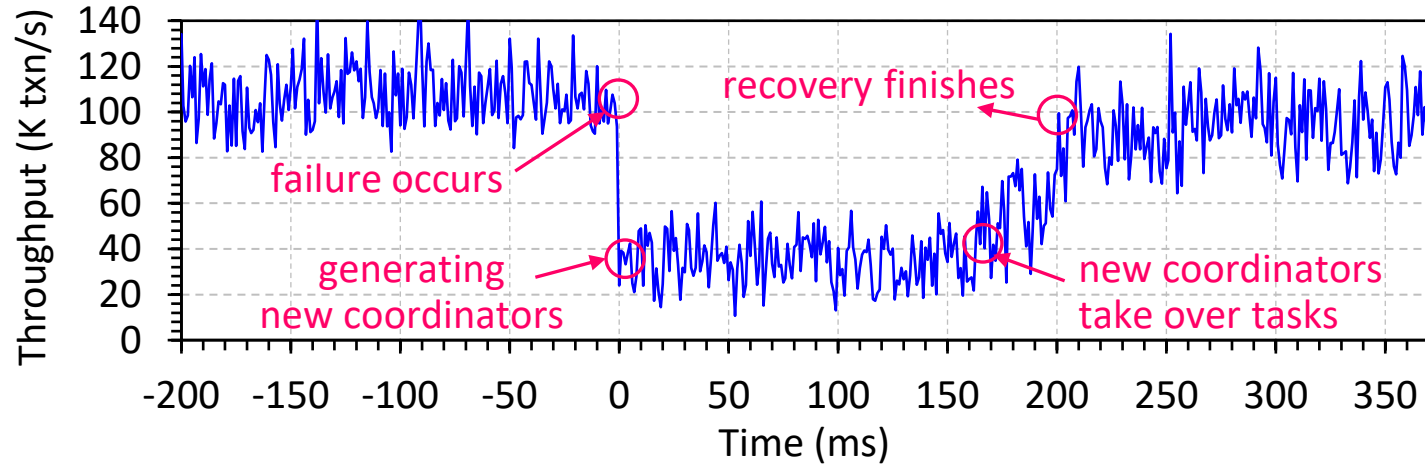


TPCC (write-intensive)

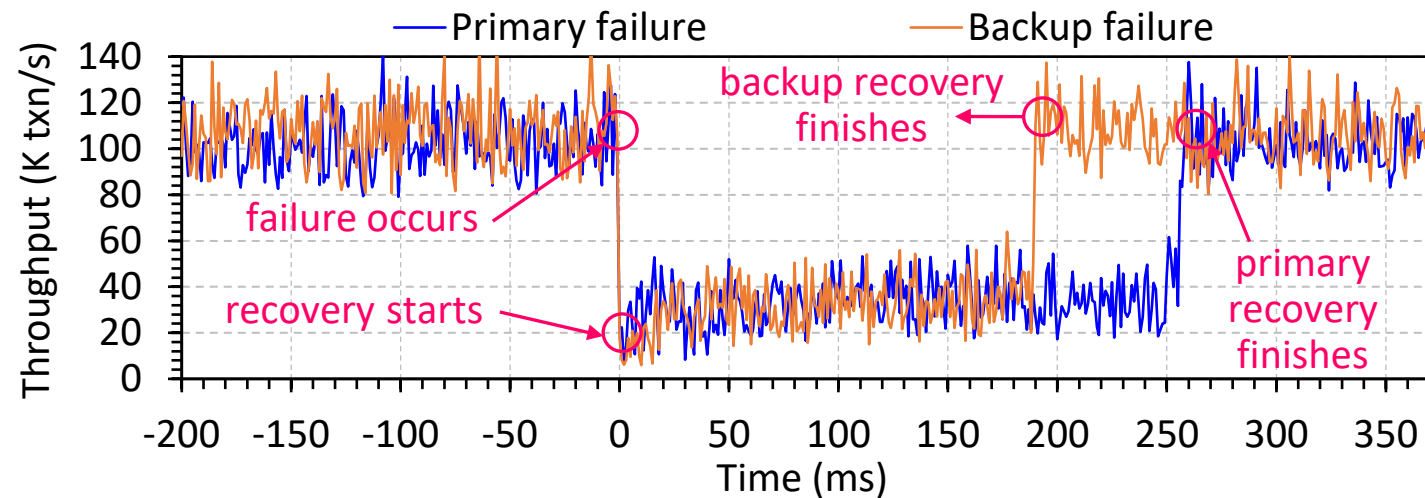


# Failure Recovery

## ➤ TPCC



Tolerating *coordinator* failures using *local operation logs*



Tolerating *replica* failures using *data migration*

# Conclusion

- Existing multi-versioning distributed transactions do not fit DM
  - Inefficient linked version chain
  - Incompatible transaction protocol
- **Motor**: a holistic multi-versioning design for DM
  - Consecutive version tuple structure (memory pool)
  - One-sided RDMA MVCC based on CVT (compute pool)
- **Benefits**

High Throughput

Low Latency

Low Memory Overhead



***Thank you! Q&A***