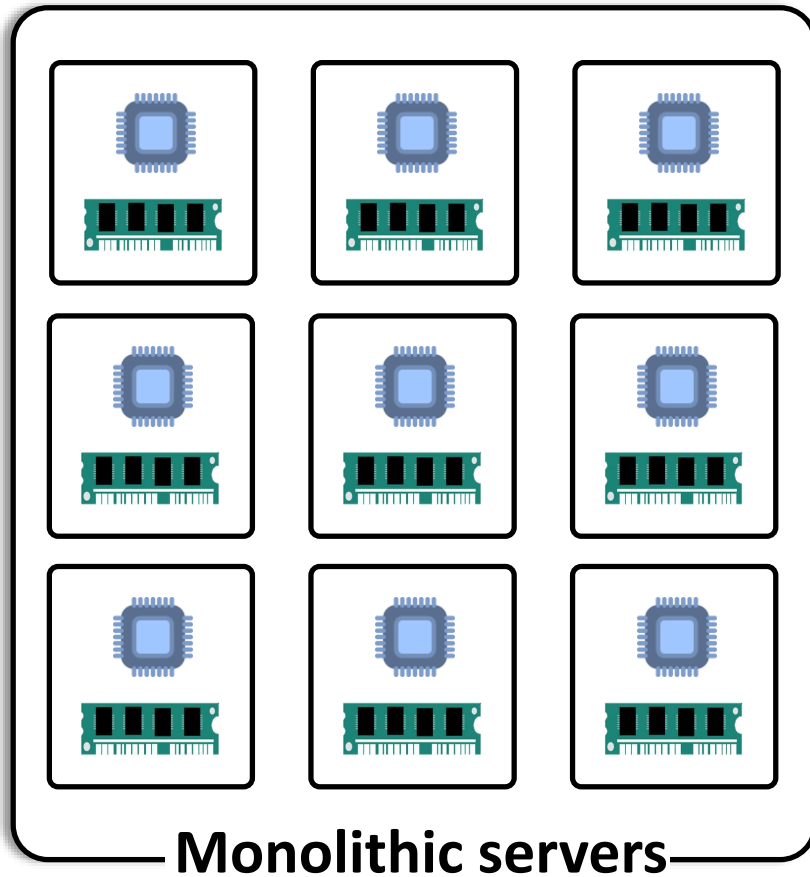


FORD: Fast One-sided RDMA-based Distributed Transactions for Disaggregated Persistent Memory

Ming Zhang, Yu Hua, Pengfei Zuo, Lurong Liu
Huazhong University of Science and Technology, China

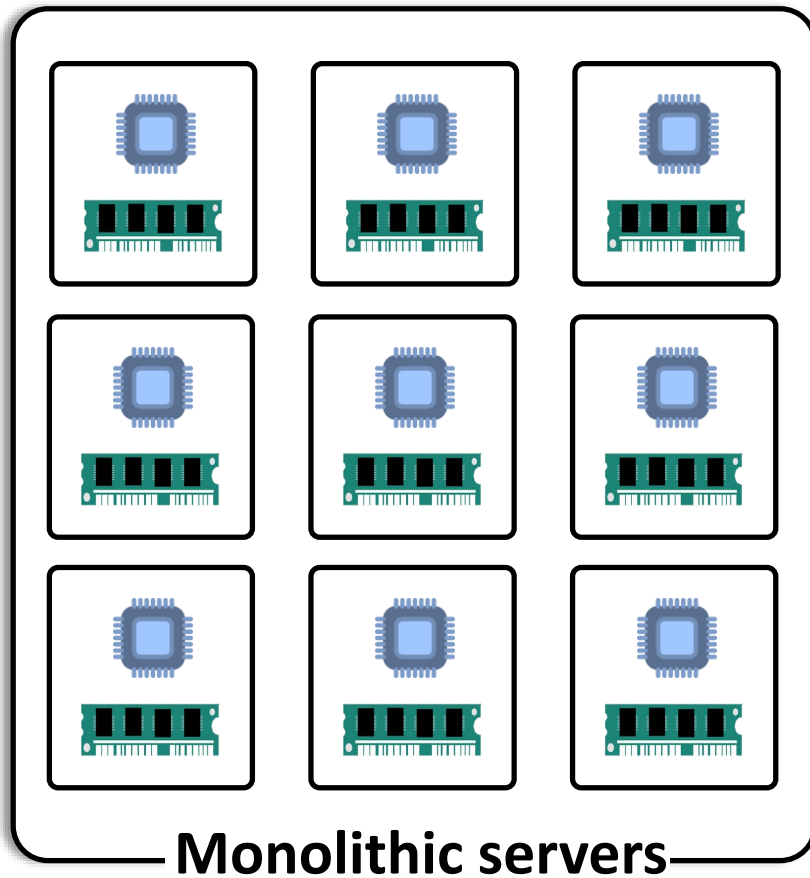
Disaggregated Persistent Memory

- Memory disaggregation

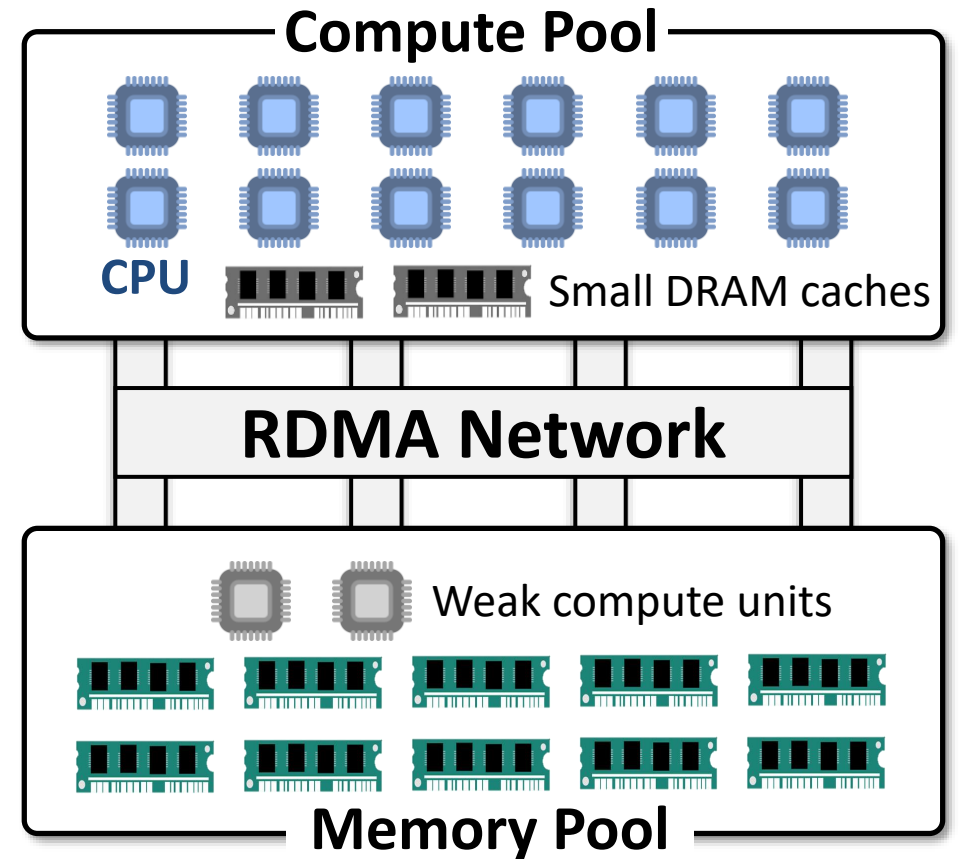


Disaggregated Persistent Memory

- Memory disaggregation

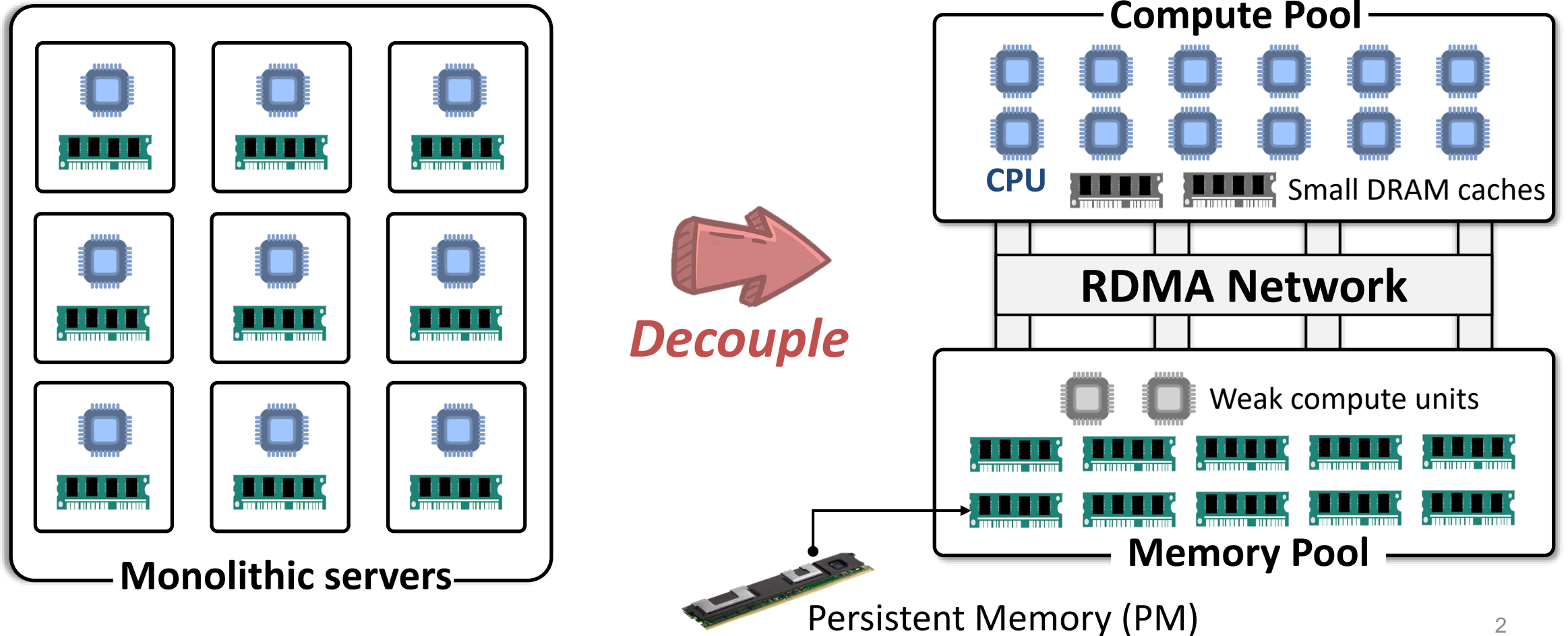


Decouple

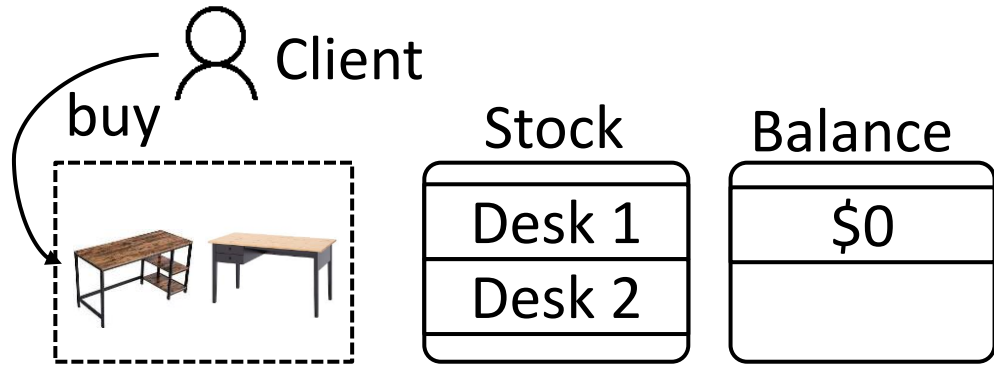


Disaggregated Persistent Memory

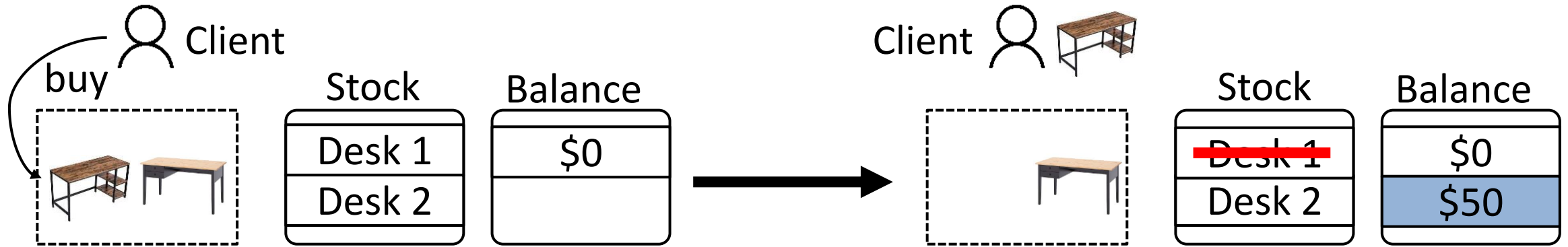
- Memory disaggregation



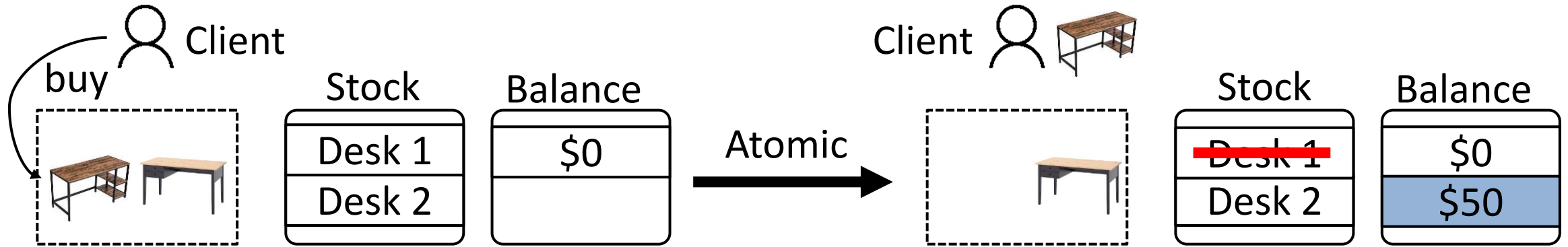
Transaction



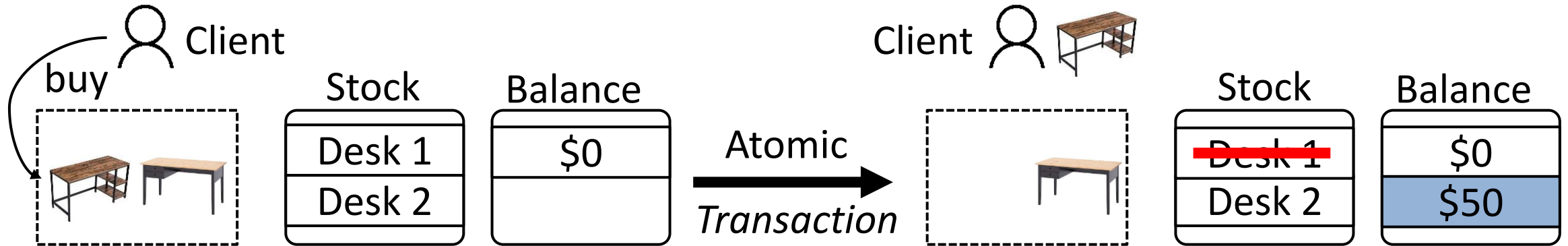
Transaction



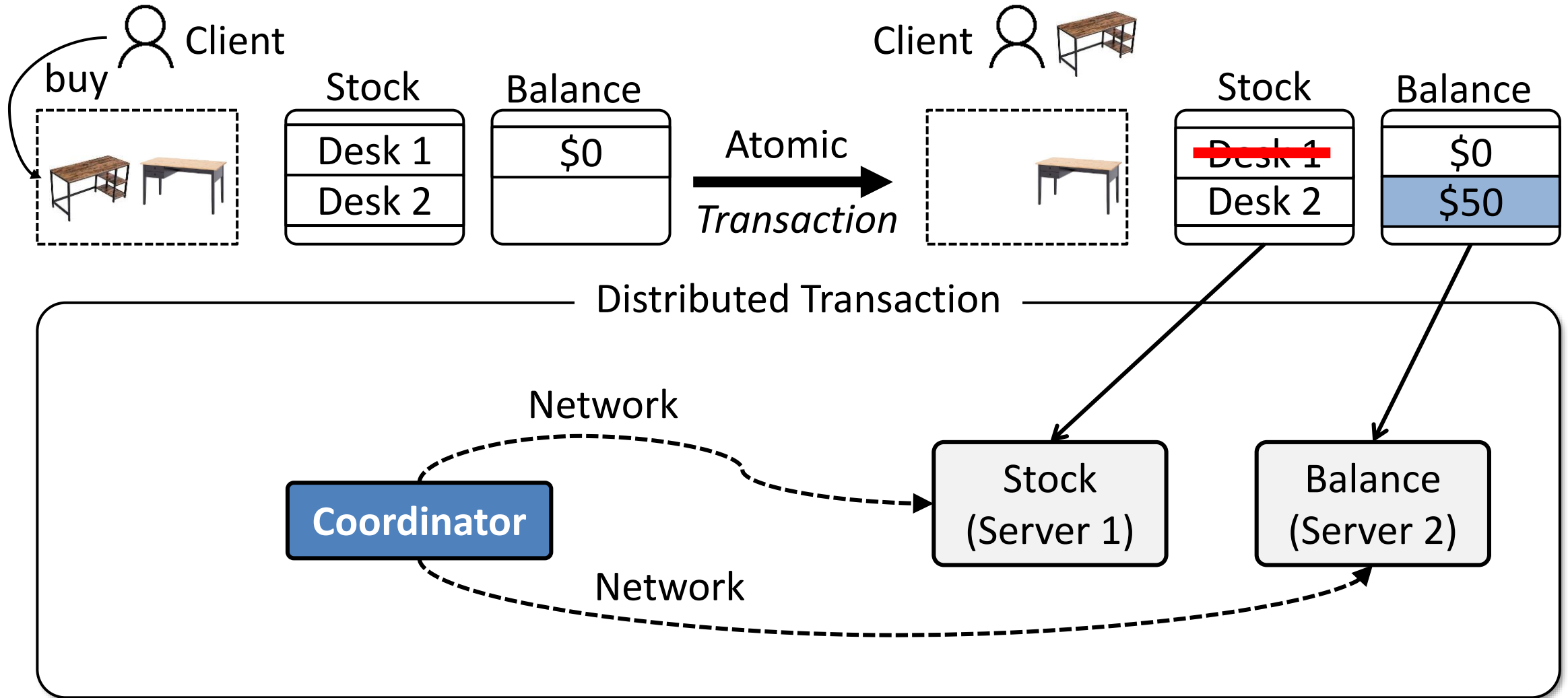
Transaction



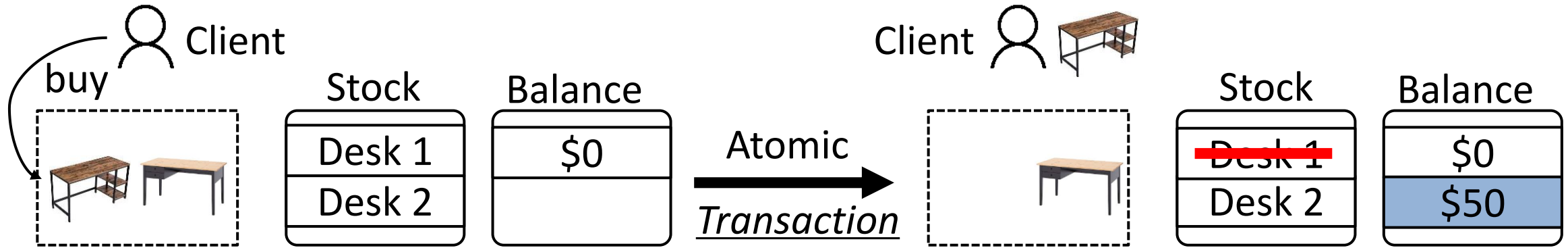
Transaction



Distributed Transaction



Distributed Transaction



Distributed Transaction: **A key building block**

Google
Cloud
Spanner



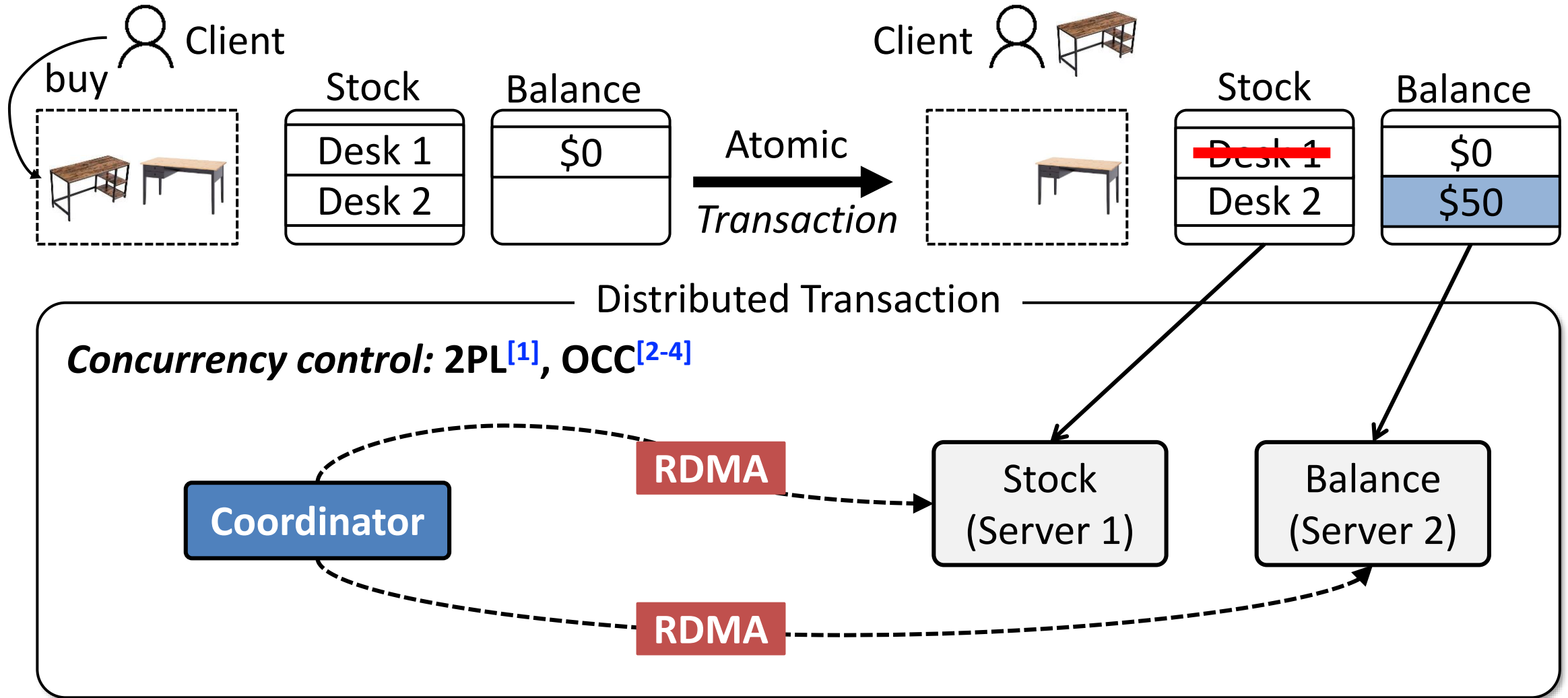
Microsoft Azure



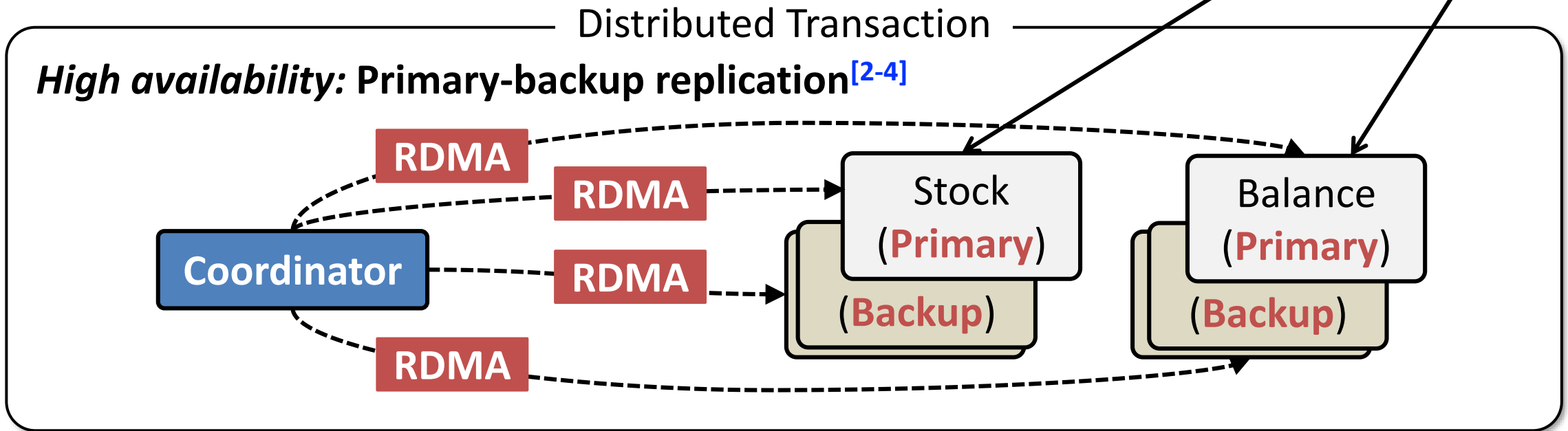
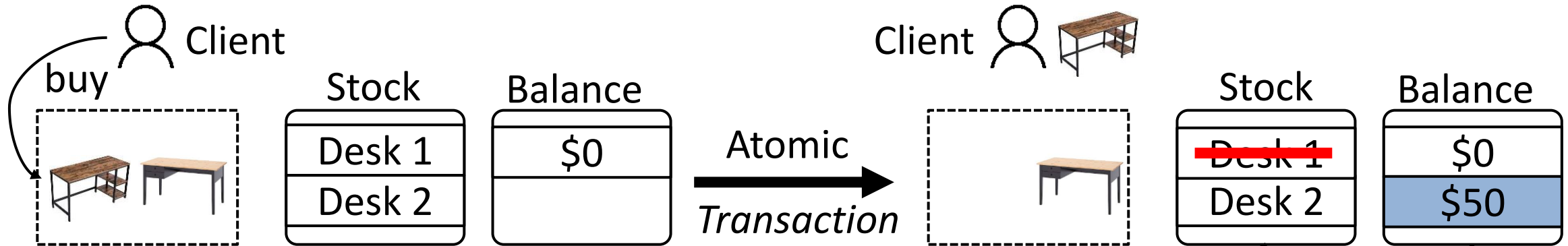
HUAWEI CLOUD



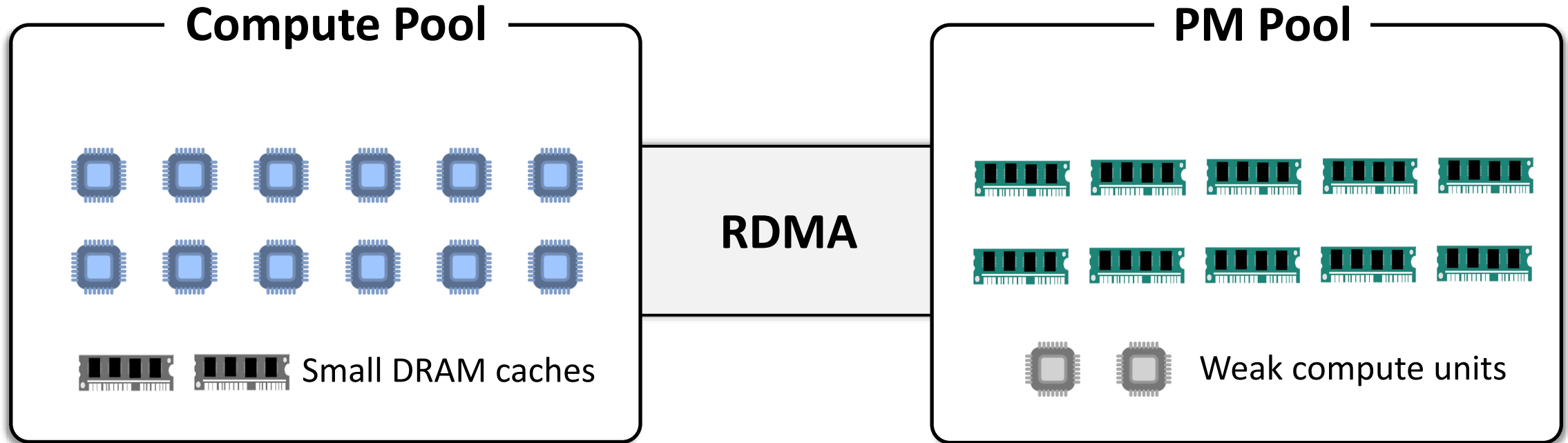
RDMA-based Distributed Transaction



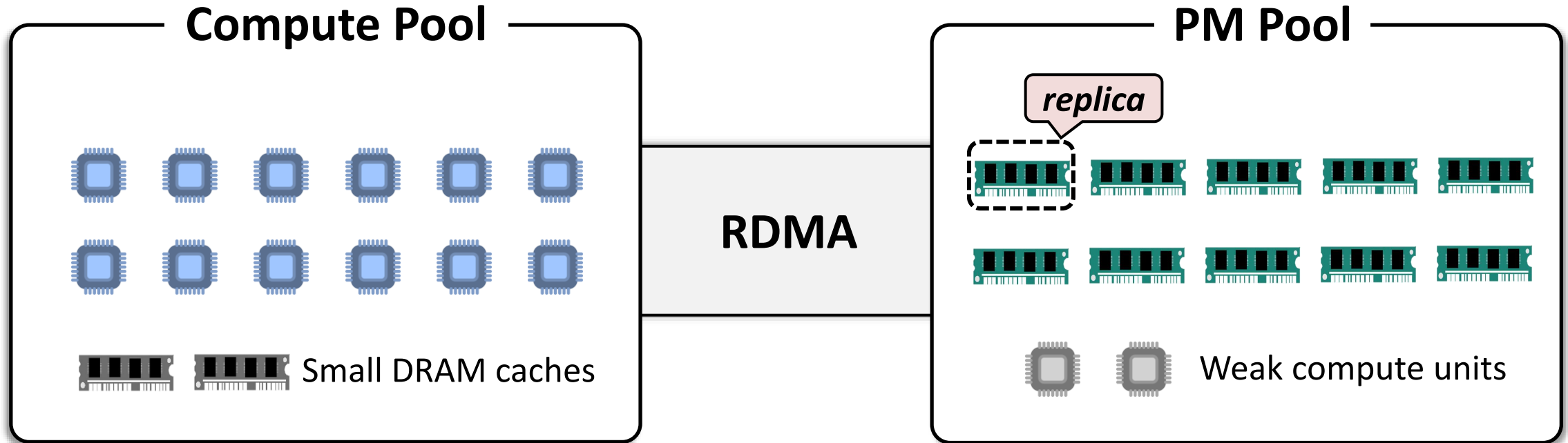
RDMA-based Distributed Transaction



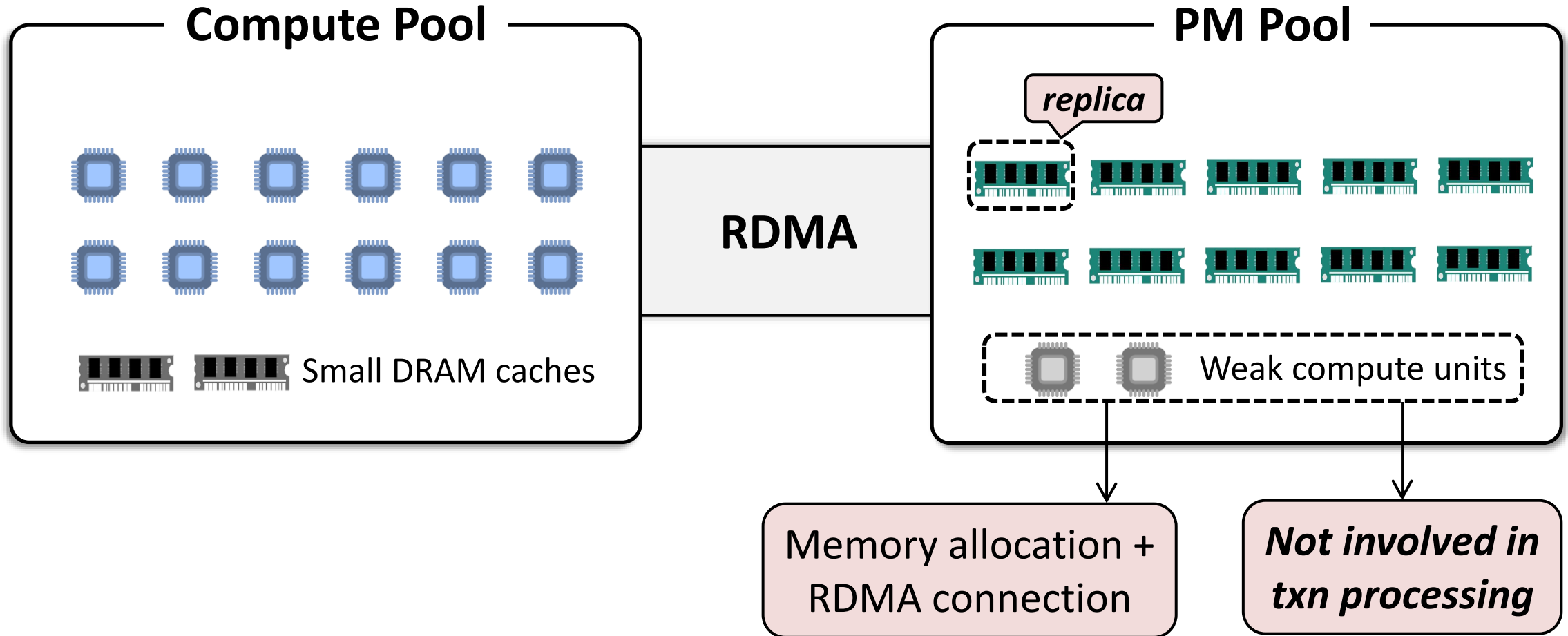
Transaction on Disaggregated PM



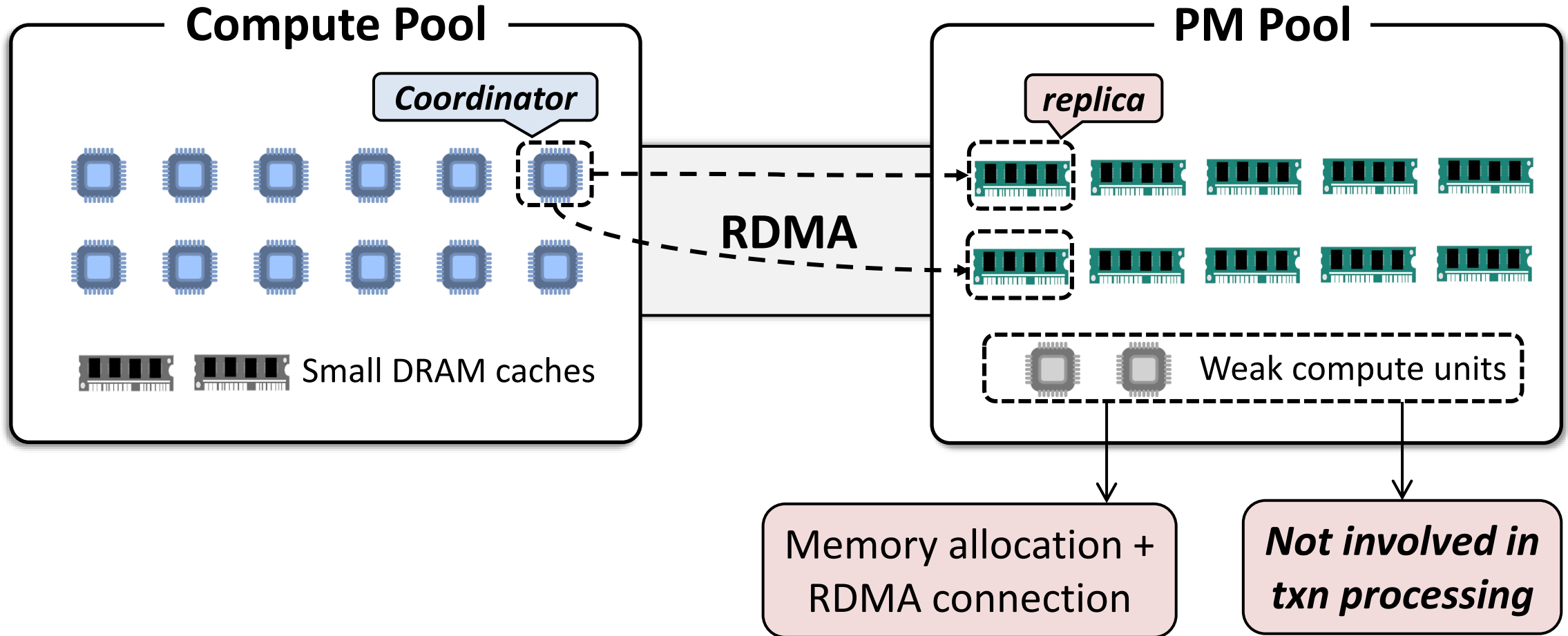
Transaction on Disaggregated PM



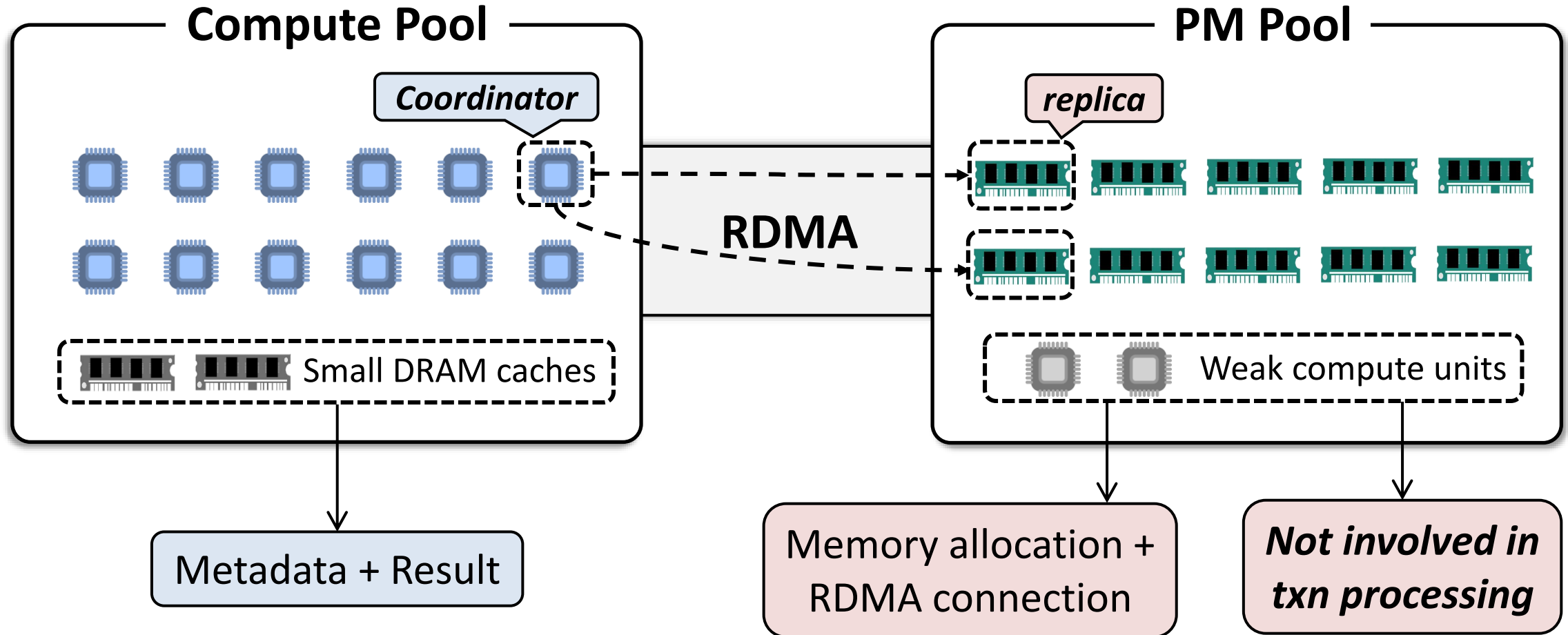
Transaction on Disaggregated PM



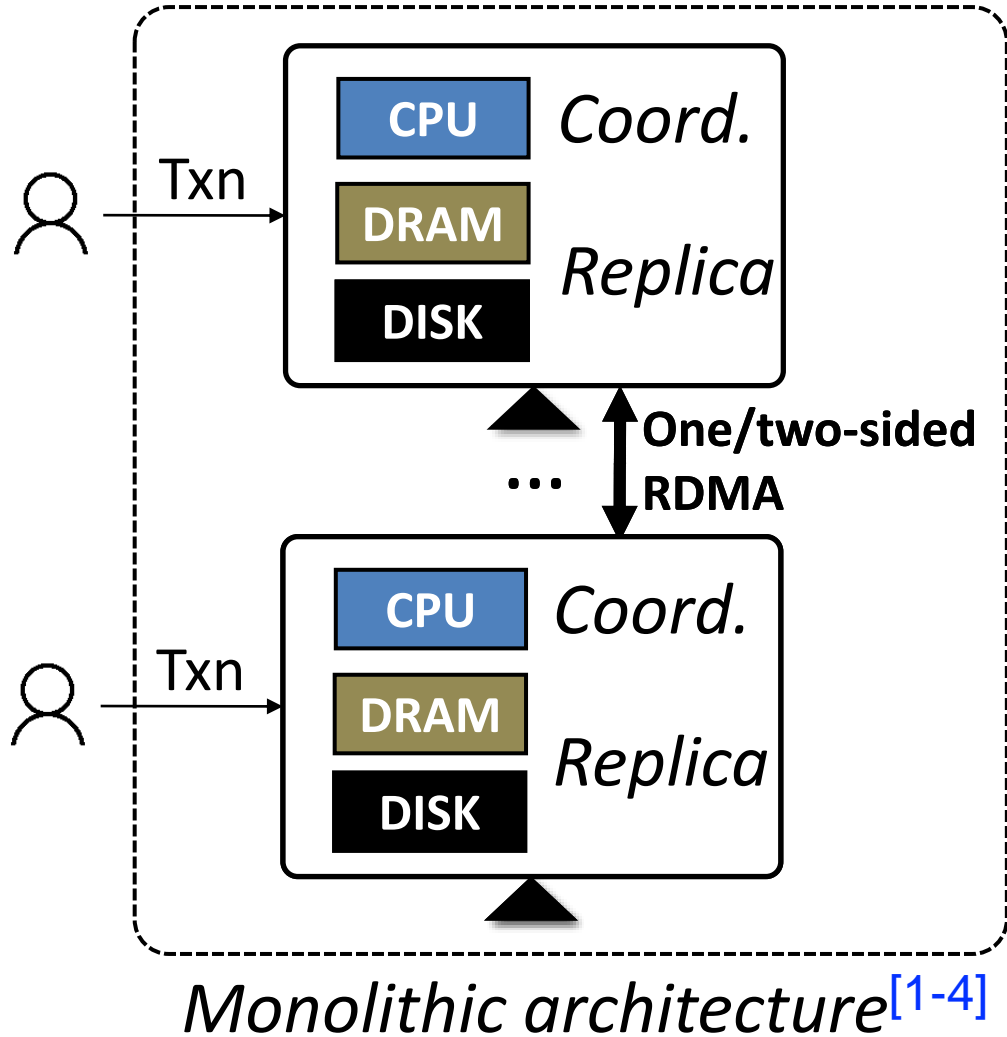
Transaction on Disaggregated PM



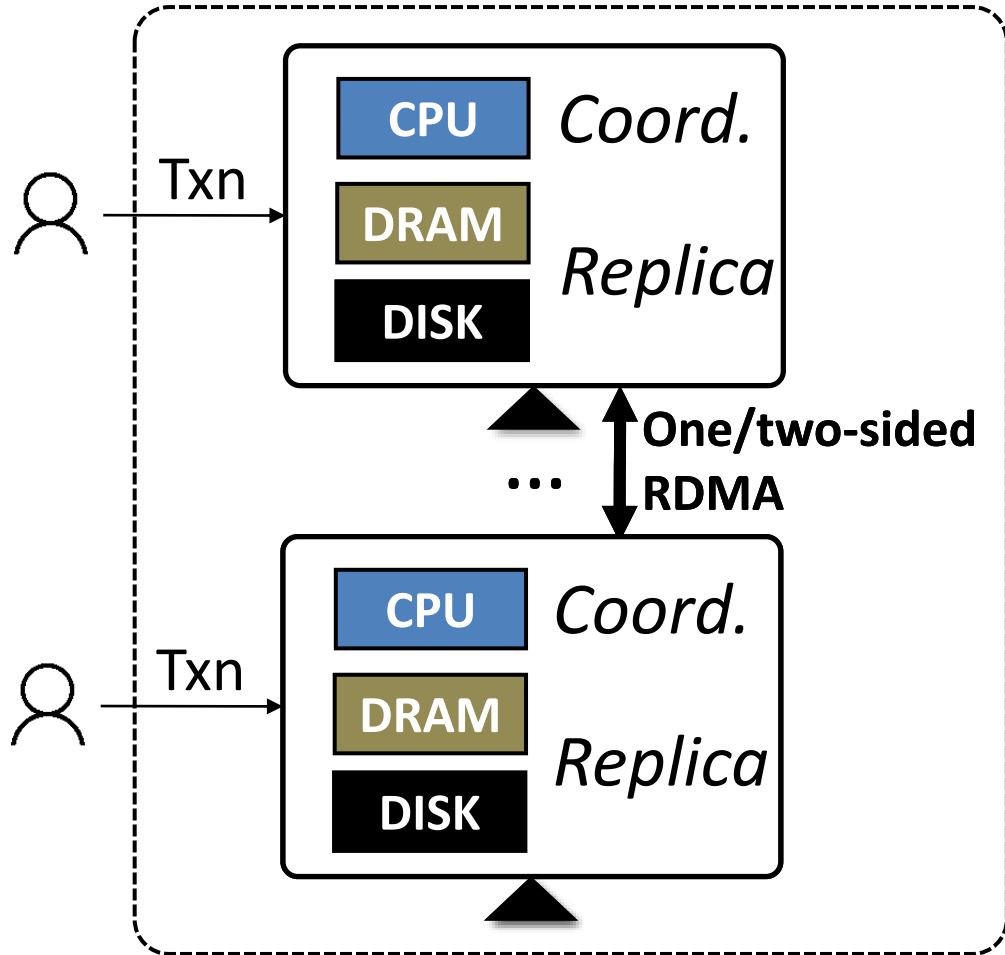
Transaction on Disaggregated PM



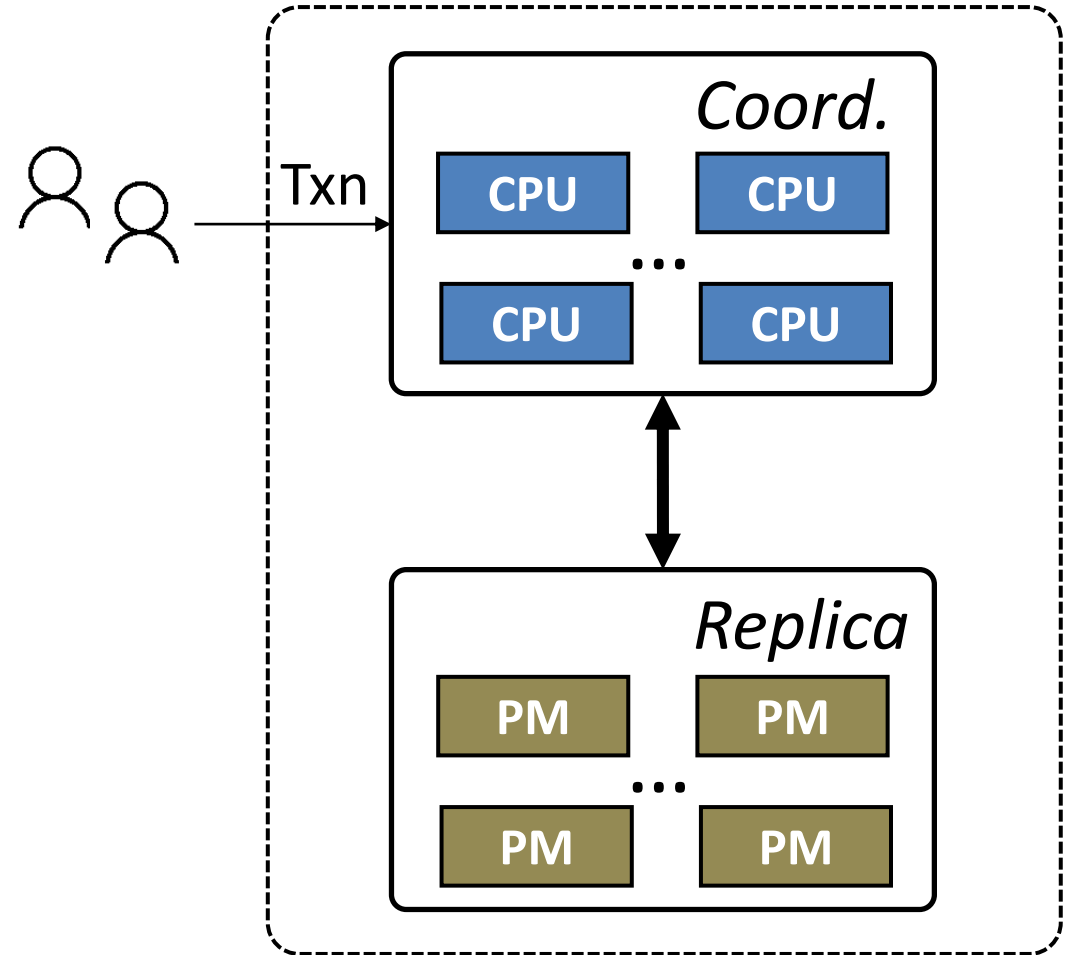
State-of-the-art



State-of-the-art

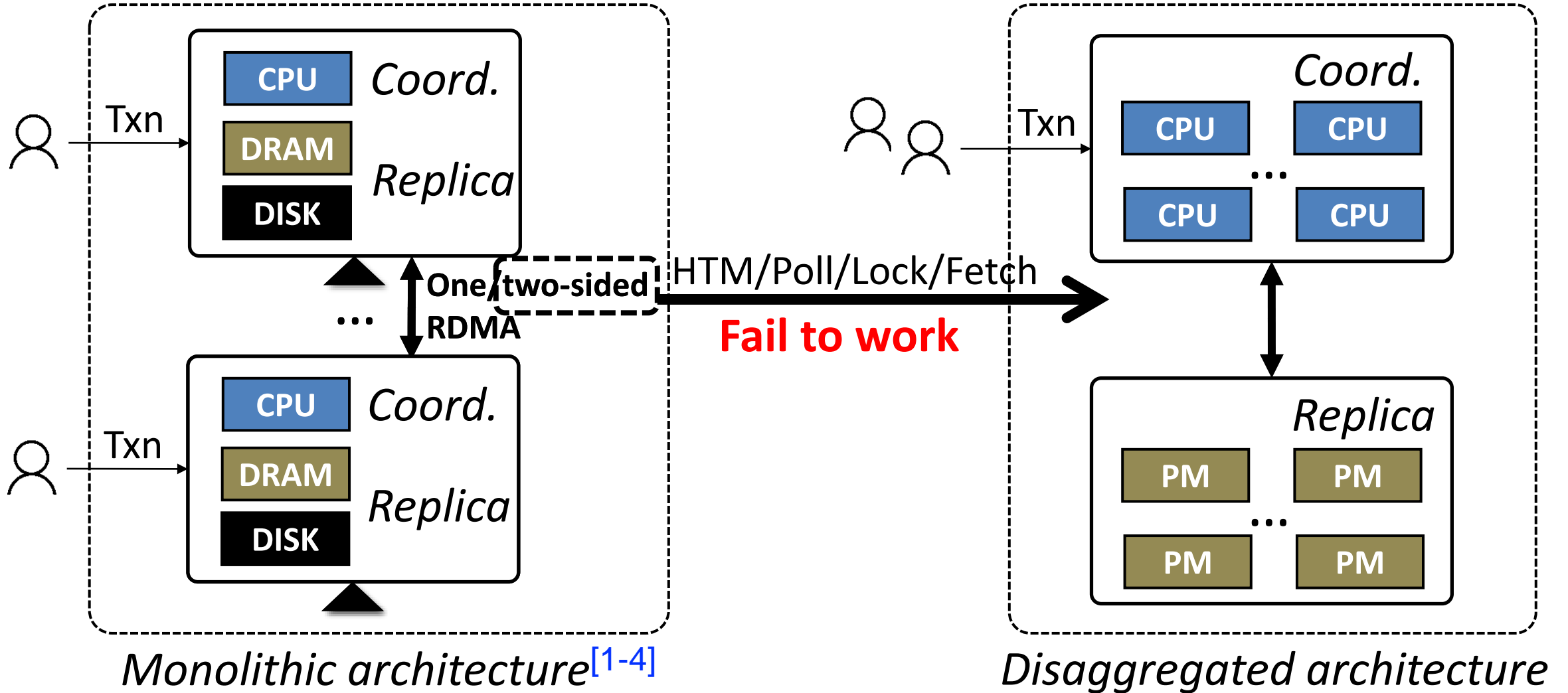


Monolithic architecture^[1-4]



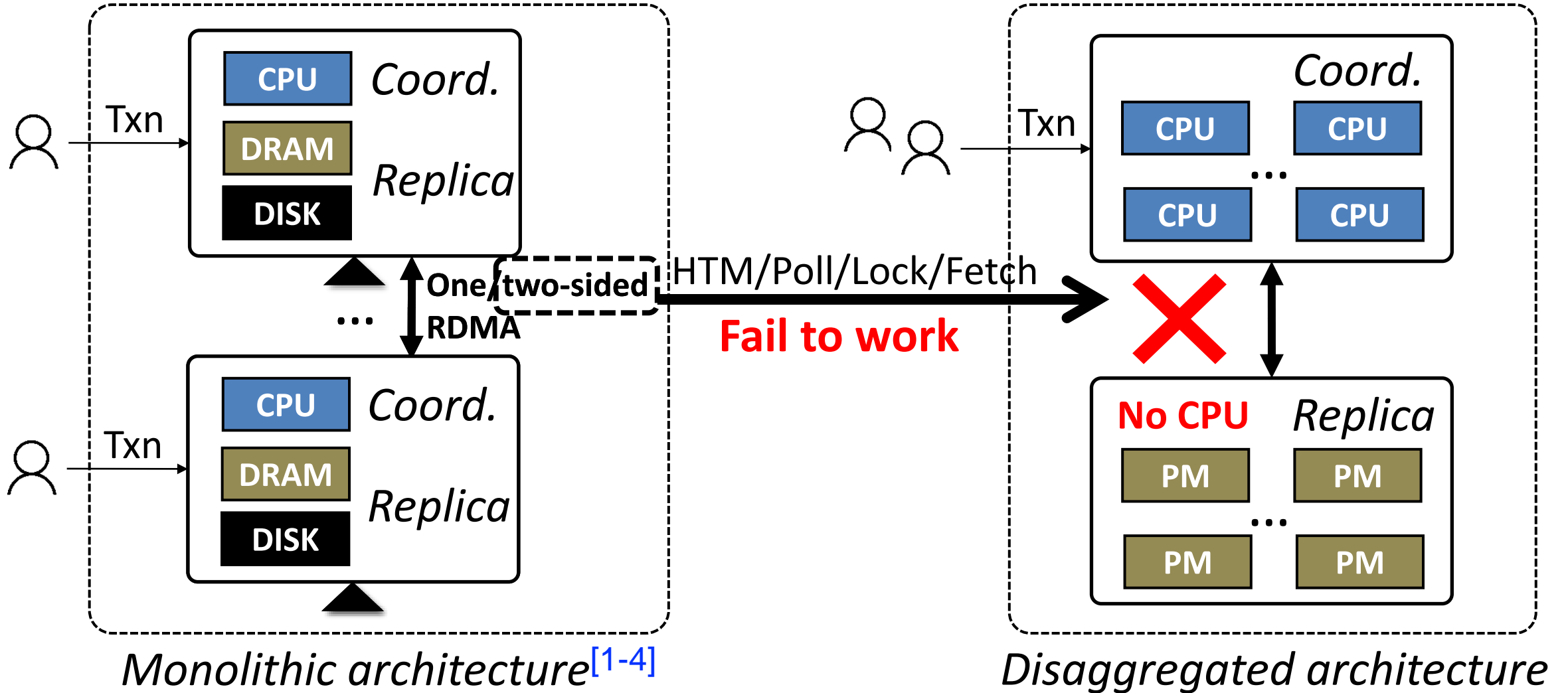
Disaggregated architecture

State-of-the-art



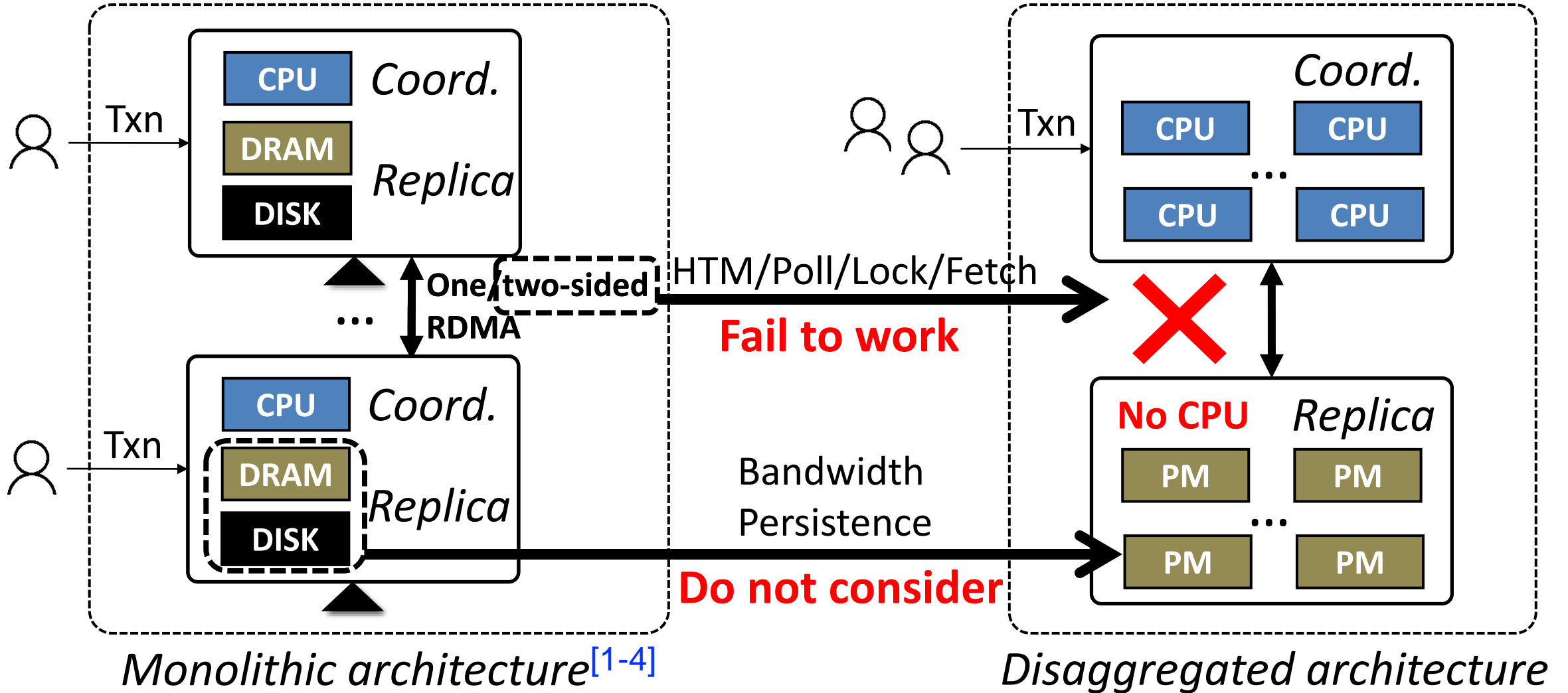
¹ DrTM@SOSP'15 ² FaRM@SOSP'15 ³ FaSST@OSDI'16 ⁴ DrTM+H@OSDI'18

State-of-the-art



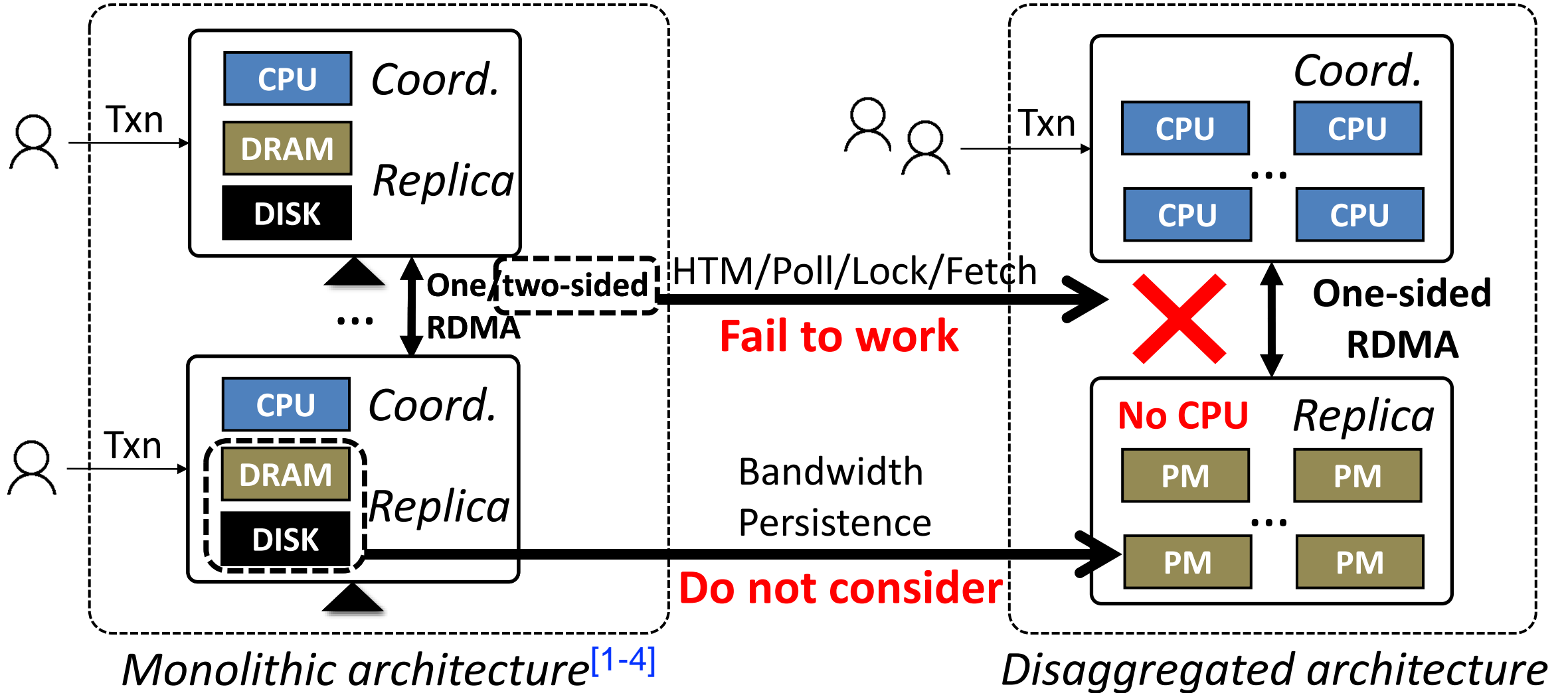
¹ DrTM@SOSP'15 ² FaRM@SOSP'15 ³ FaSST@OSDI'16 ⁴ DrTM+H@OSDI'18

State-of-the-art



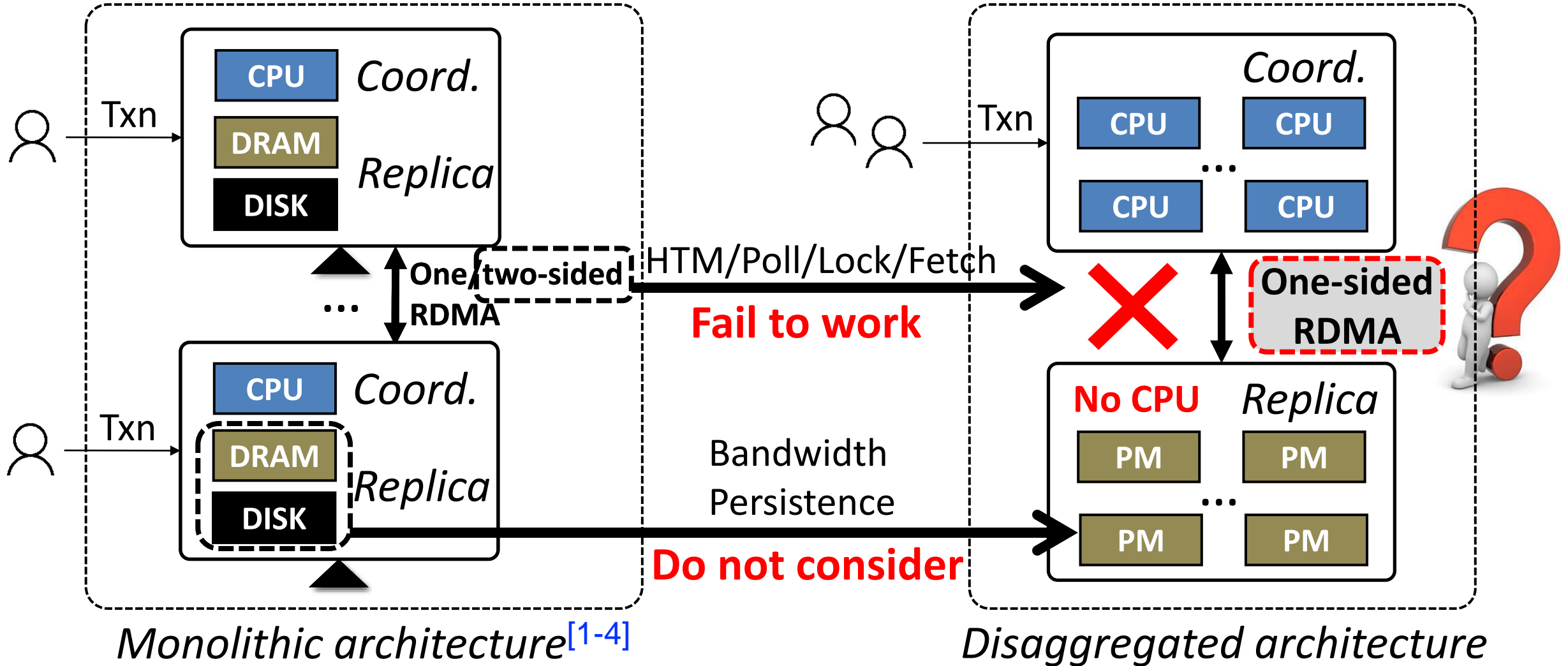
¹ DrTM@SOSP'15 ² FaRM@SOSP'15 ³ FaSST@OSDI'16 ⁴ DrTM+H@OSDI'18

State-of-the-art



¹ DrTM@SOSP'15 ² FaRM@SOSP'15 ³ FaSST@OSDI'16 ⁴ DrTM+H@OSDI'18

State-of-the-art



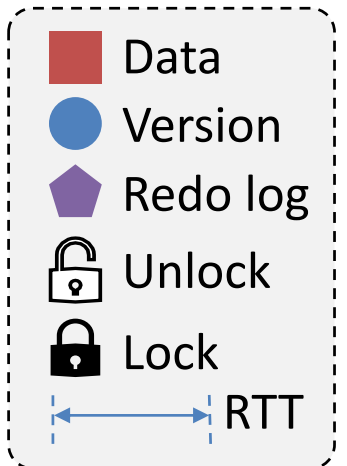
¹ DrTM@SOSP'15 ² FaRM@SOSP'15 ³ FaSST@OSDI'16 ⁴ DrTM+H@OSDI'18

Challenge 1

➤ Long-latency processing: *Many round trips*

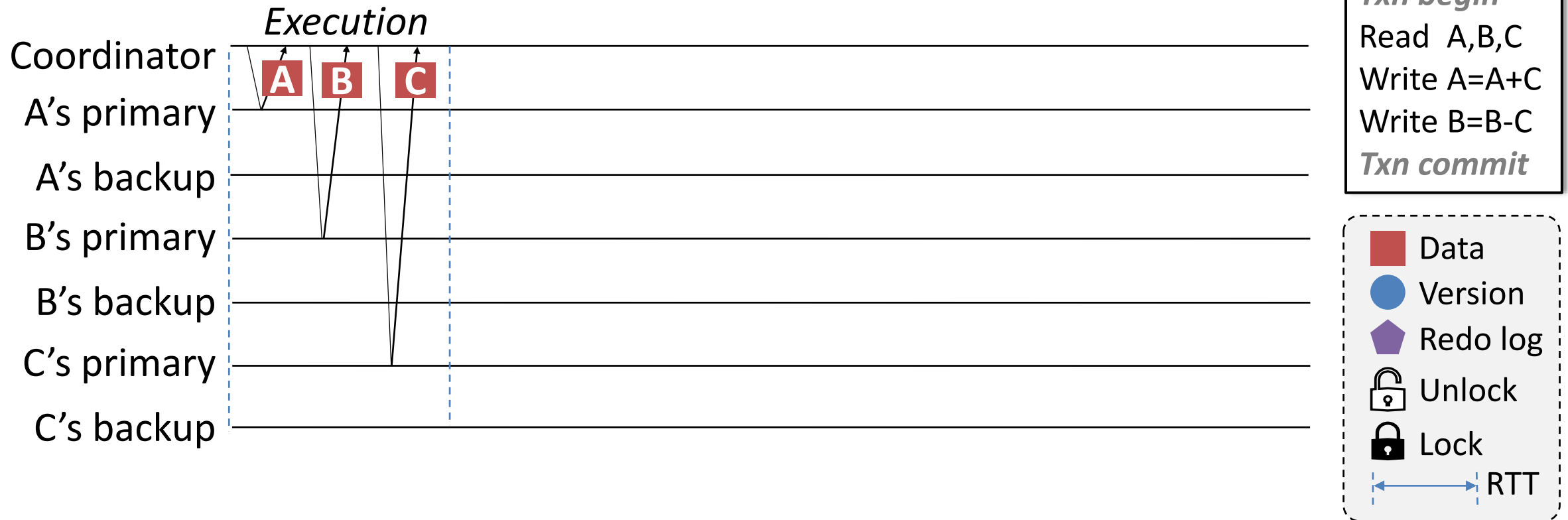


Txn begin
Read A,B,C
Write A=A+C
Write B=B-C
Txn commit



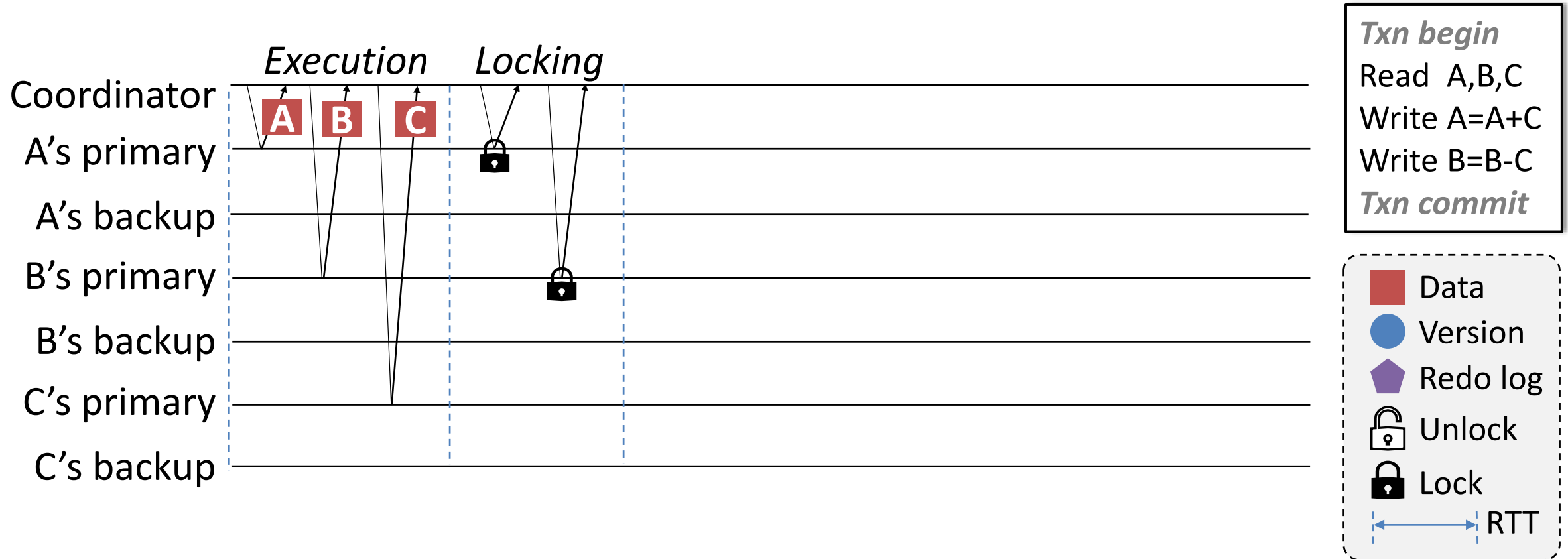
Challenge 1

➤ Long-latency processing: *Many round trips*



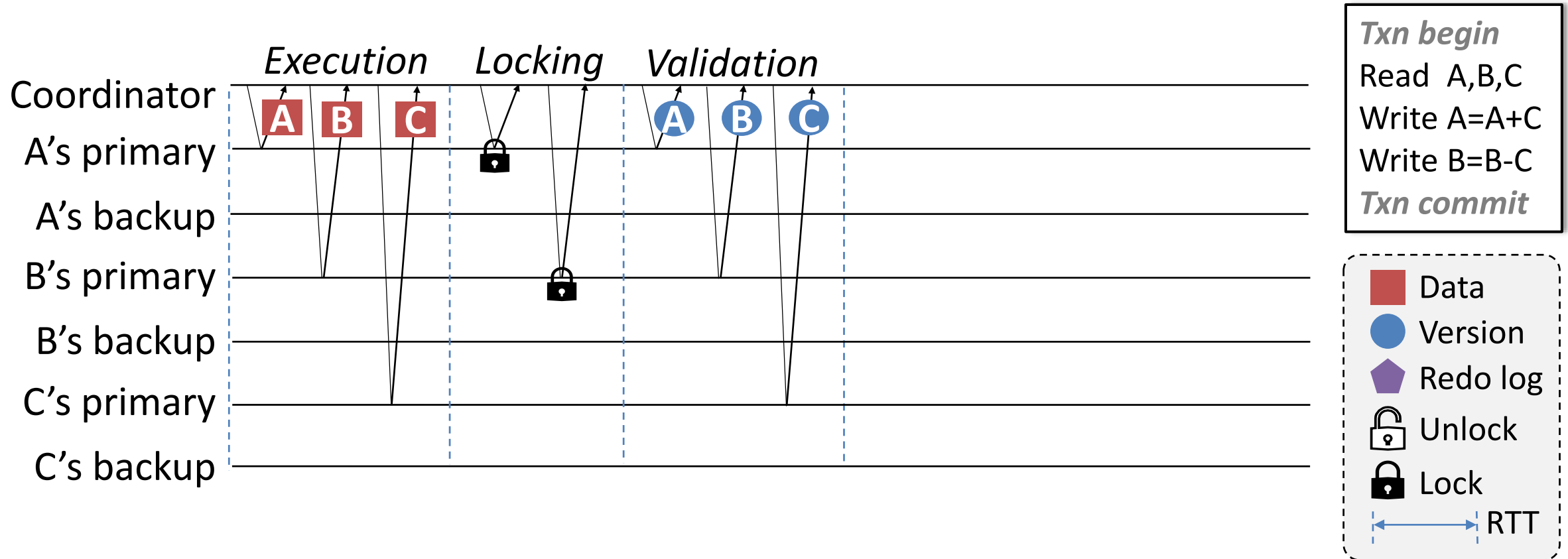
Challenge 1

➤ Long-latency processing: *Many round trips*



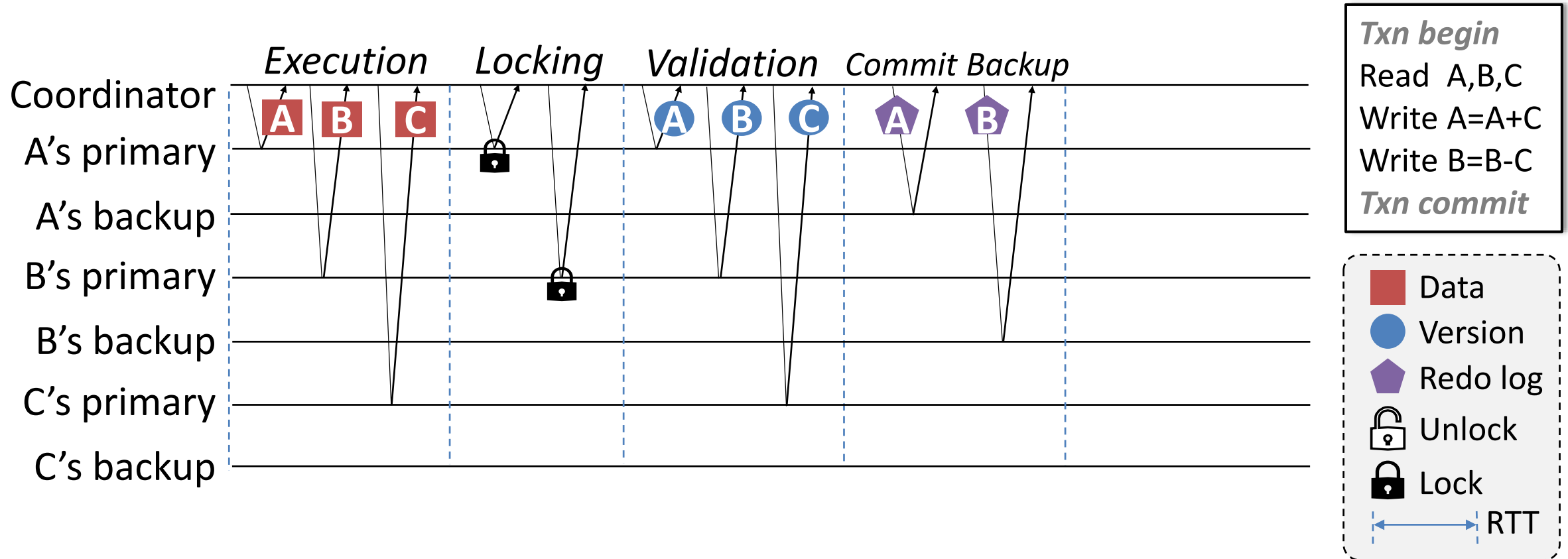
Challenge 1

➤ Long-latency processing: *Many round trips*



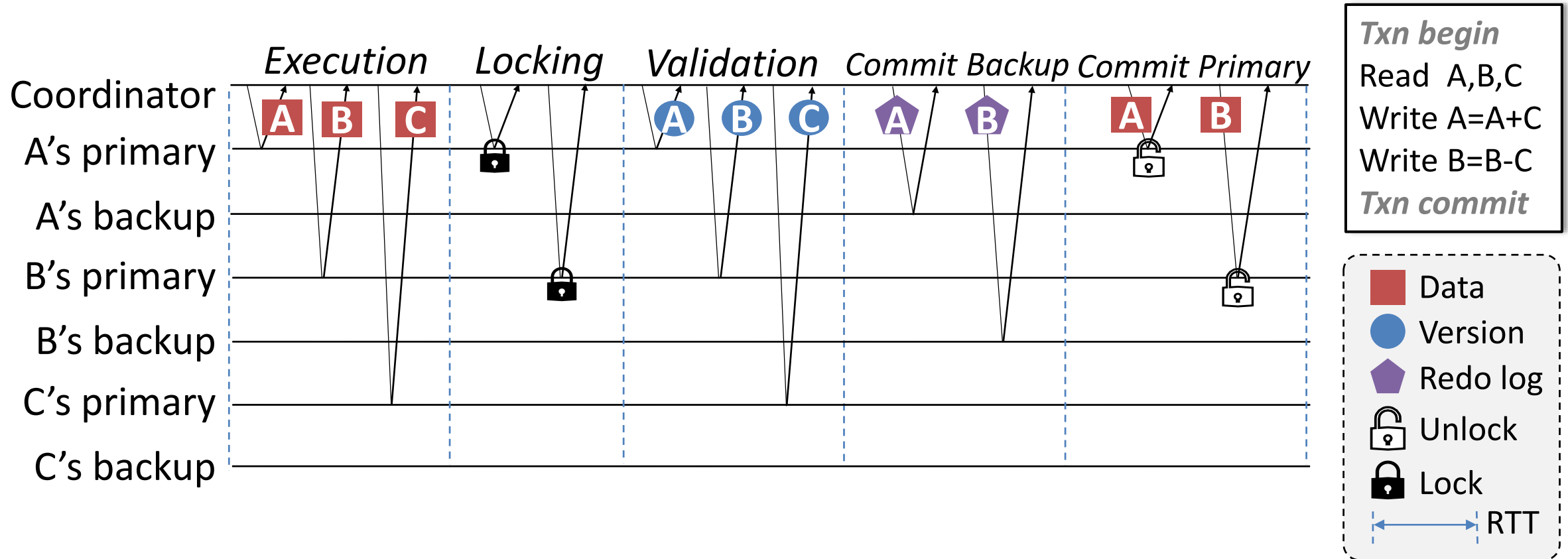
Challenge 1

➤ Long-latency processing: *Many round trips*



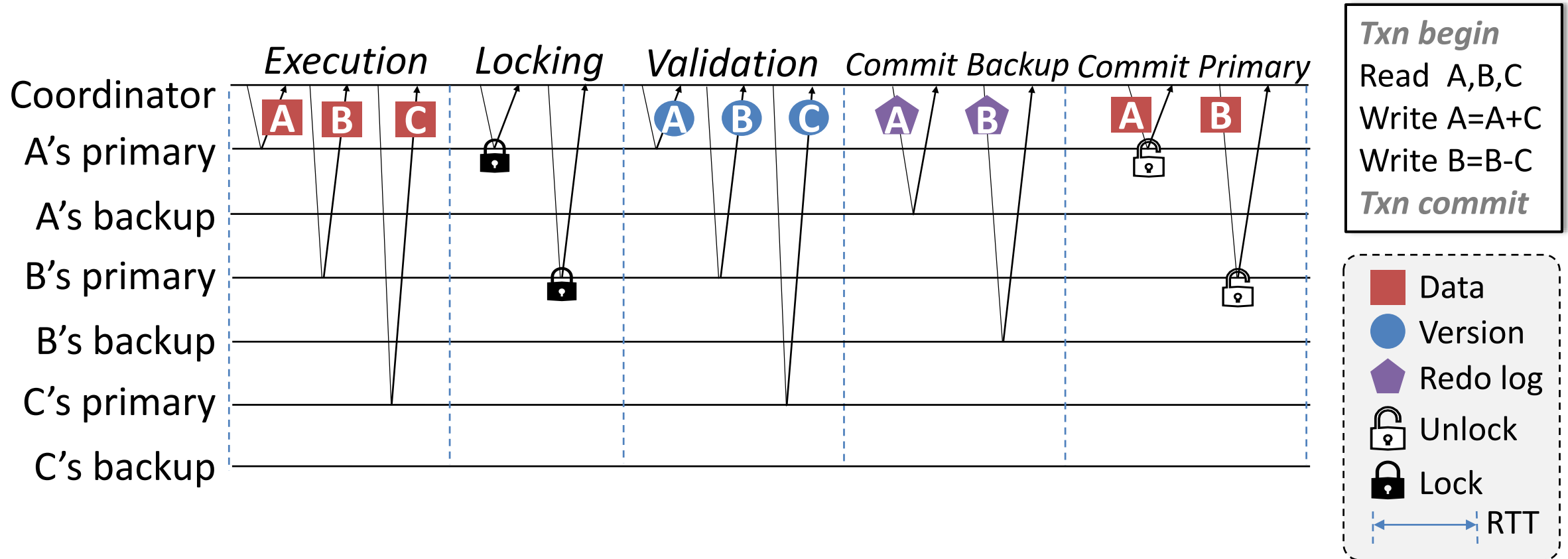
Challenge 1

➤ Long-latency processing: *Many round trips*



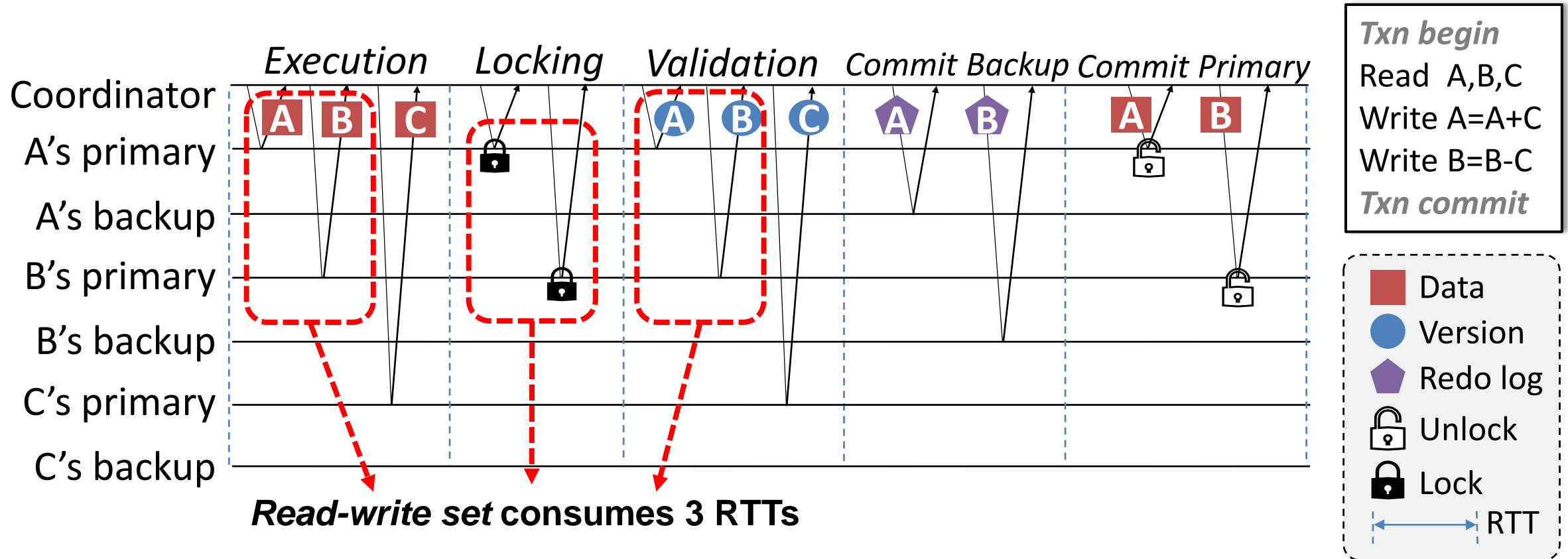
Challenge 1

➤ Long-latency processing: *Many round trips*



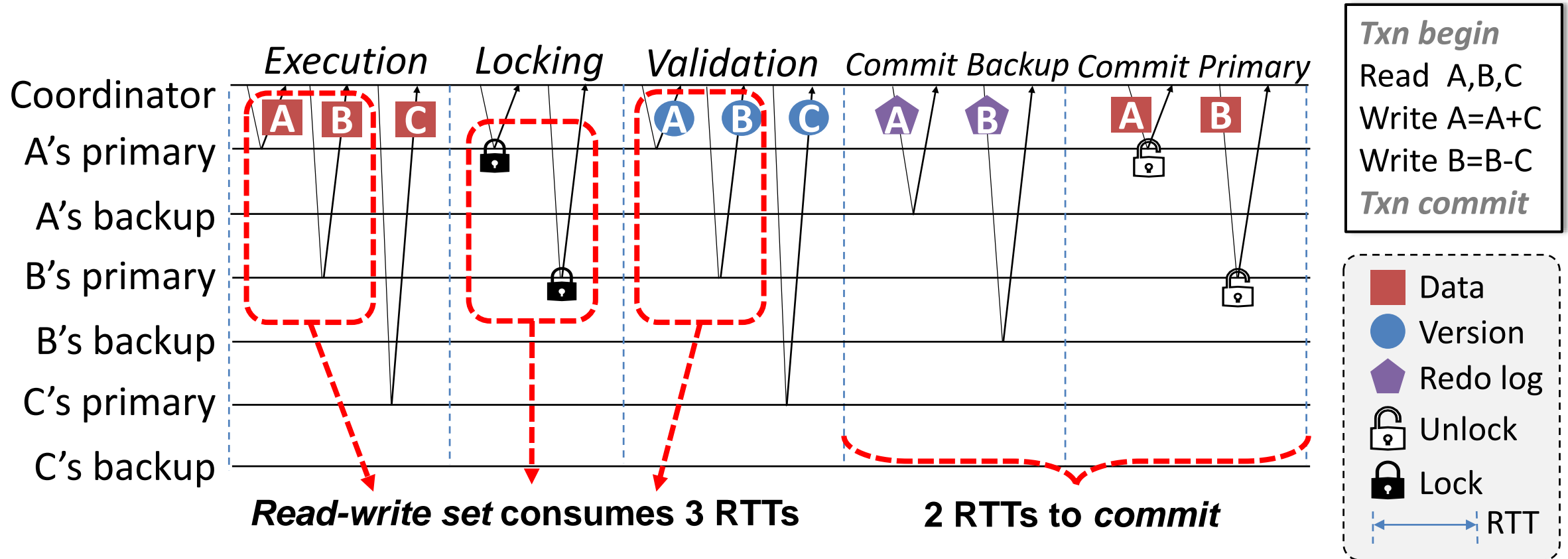
Challenge 1

➤ Long-latency processing: *Many round trips*



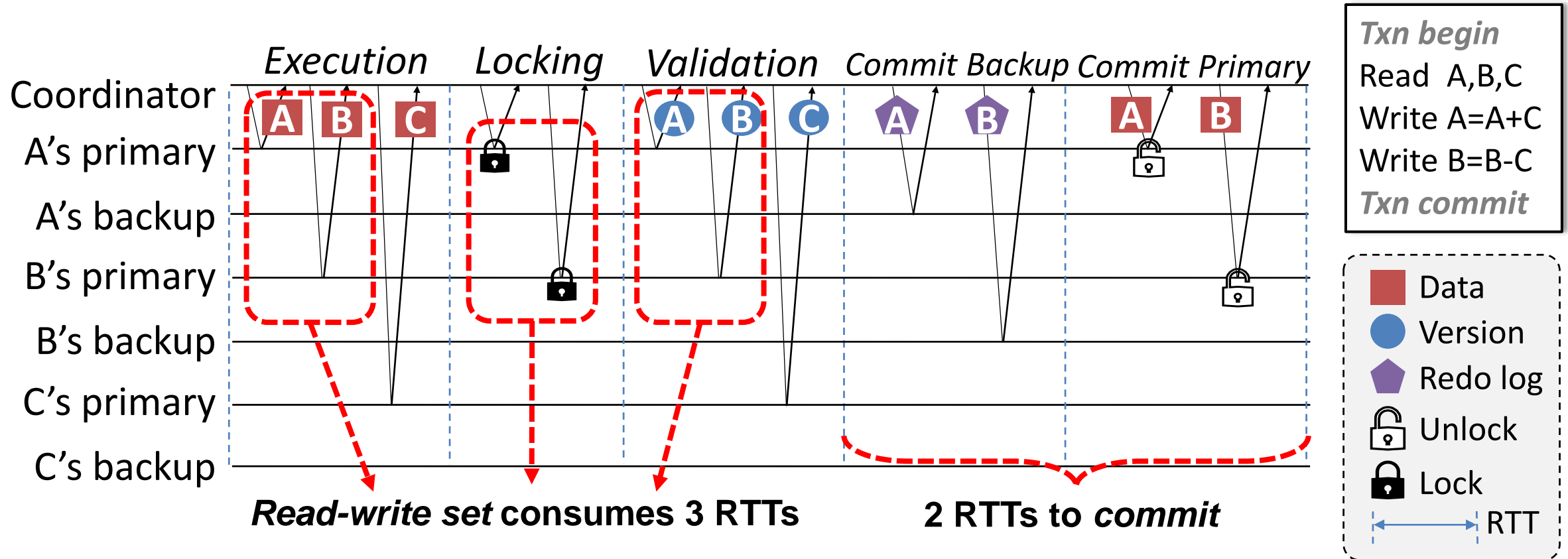
Challenge 1

➤ Long-latency processing: *Many round trips*



Challenge 1

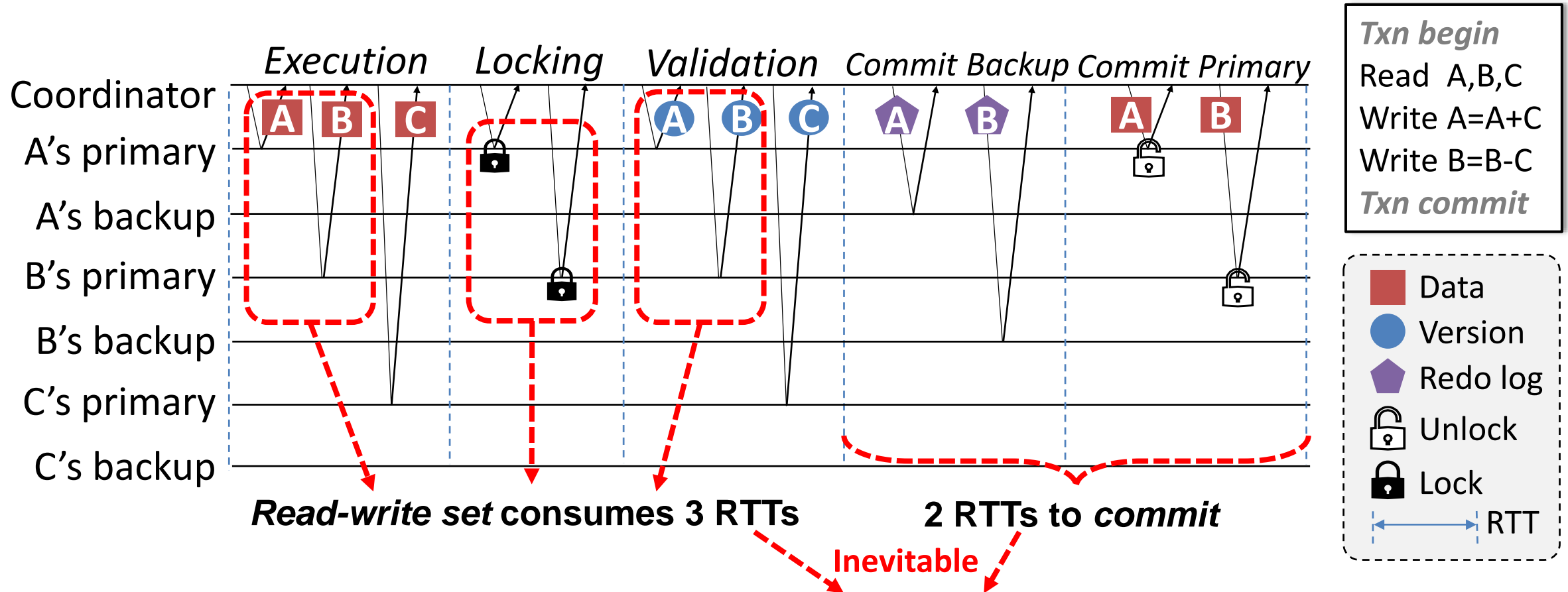
➤ Long-latency processing: *Many round trips*



Coordinators and replicas are separated: **All transactions are distributed**

Challenge 1

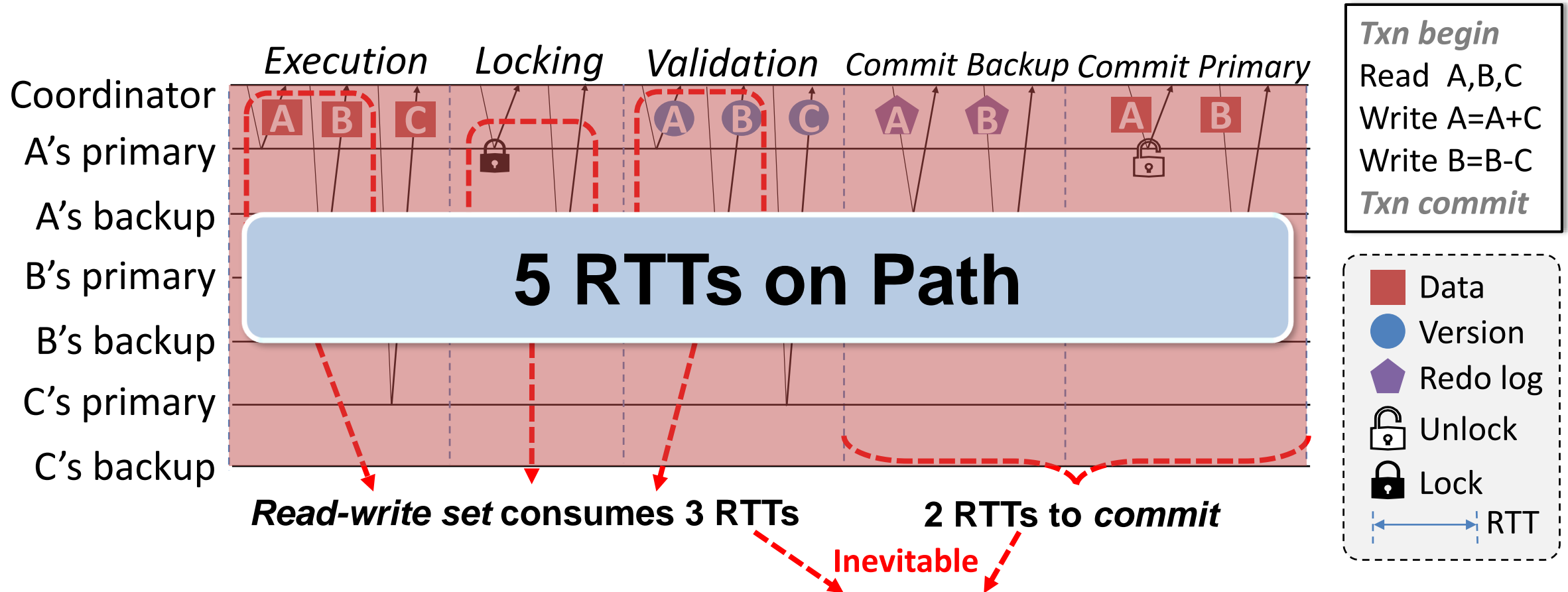
➤ Long-latency processing: *Many round trips*



Coordinators and replicas are separated: **All transactions are distributed**

Challenge 1

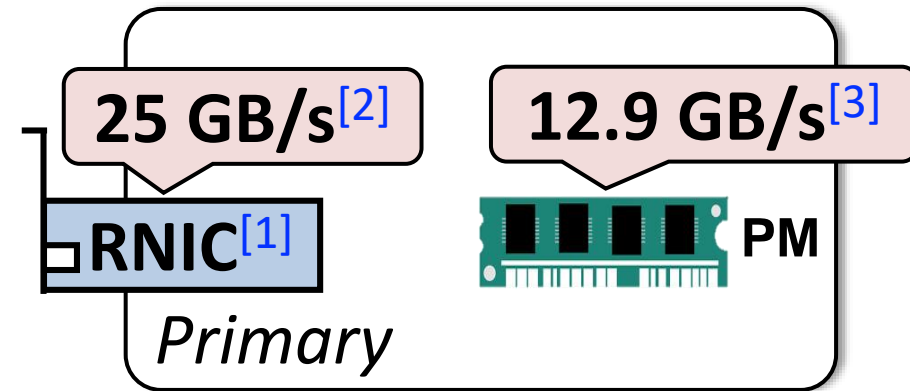
➤ Long-latency processing: *Many round trips*



Coordinators and replicas are separated: **All transactions are distributed**

Challenge 2

- Limited PM bandwidth: *High loads on primary*

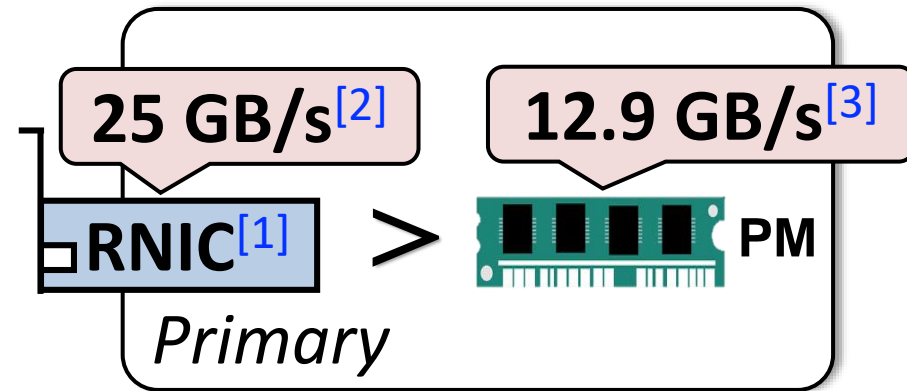


¹ RDMA NIC ² Mellanox ConnectX-6/dual-port ConnectX-5 VPI Card

³ Six interleaved 256GB Intel Optane PM [Jian Yang, et al.@FAST'20]

Challenge 2

- Limited PM bandwidth: *High loads on primary*

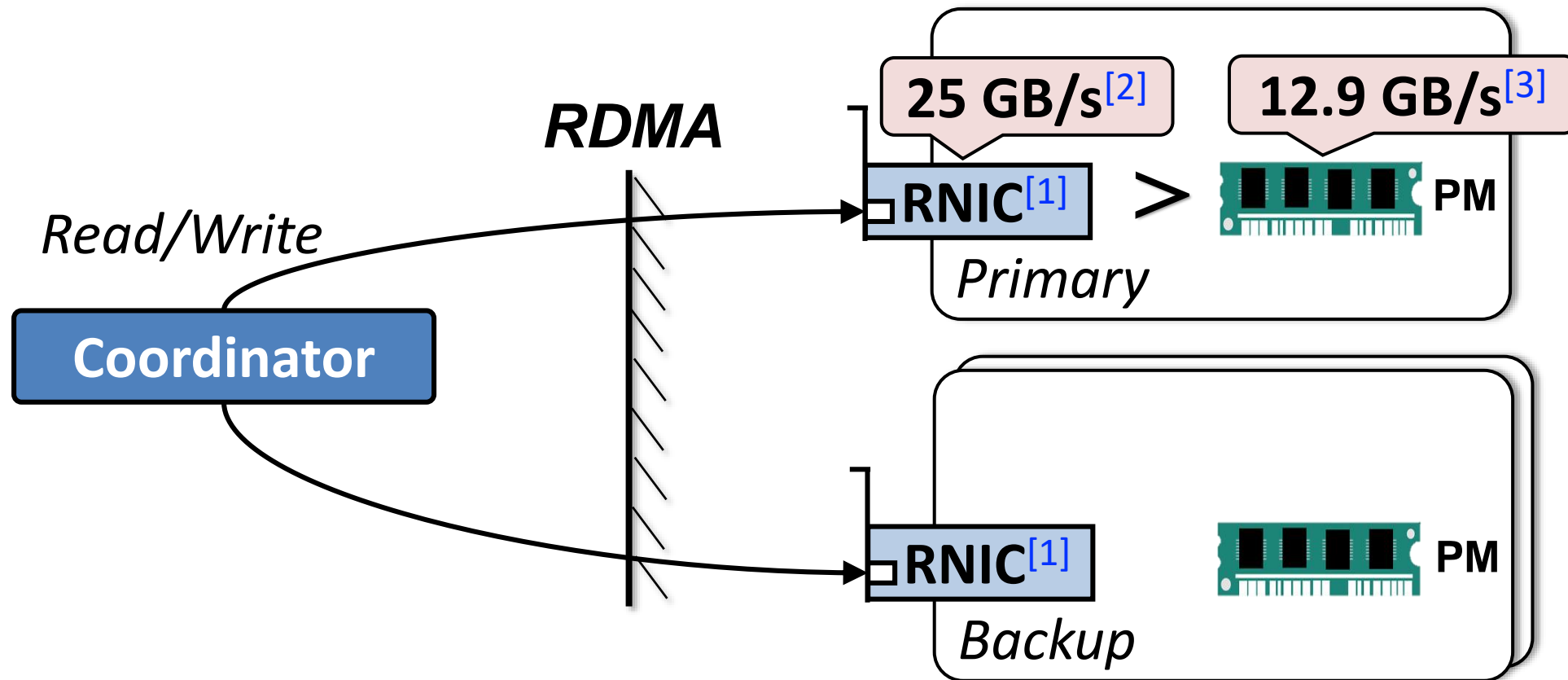


¹ RDMA NIC ² Mellanox ConnectX-6/dual-port ConnectX-5 VPI Card

³ Six interleaved 256GB Intel Optane PM [Jian Yang, et al.@FAST'20]

Challenge 2

- Limited PM bandwidth: *High loads on primary*

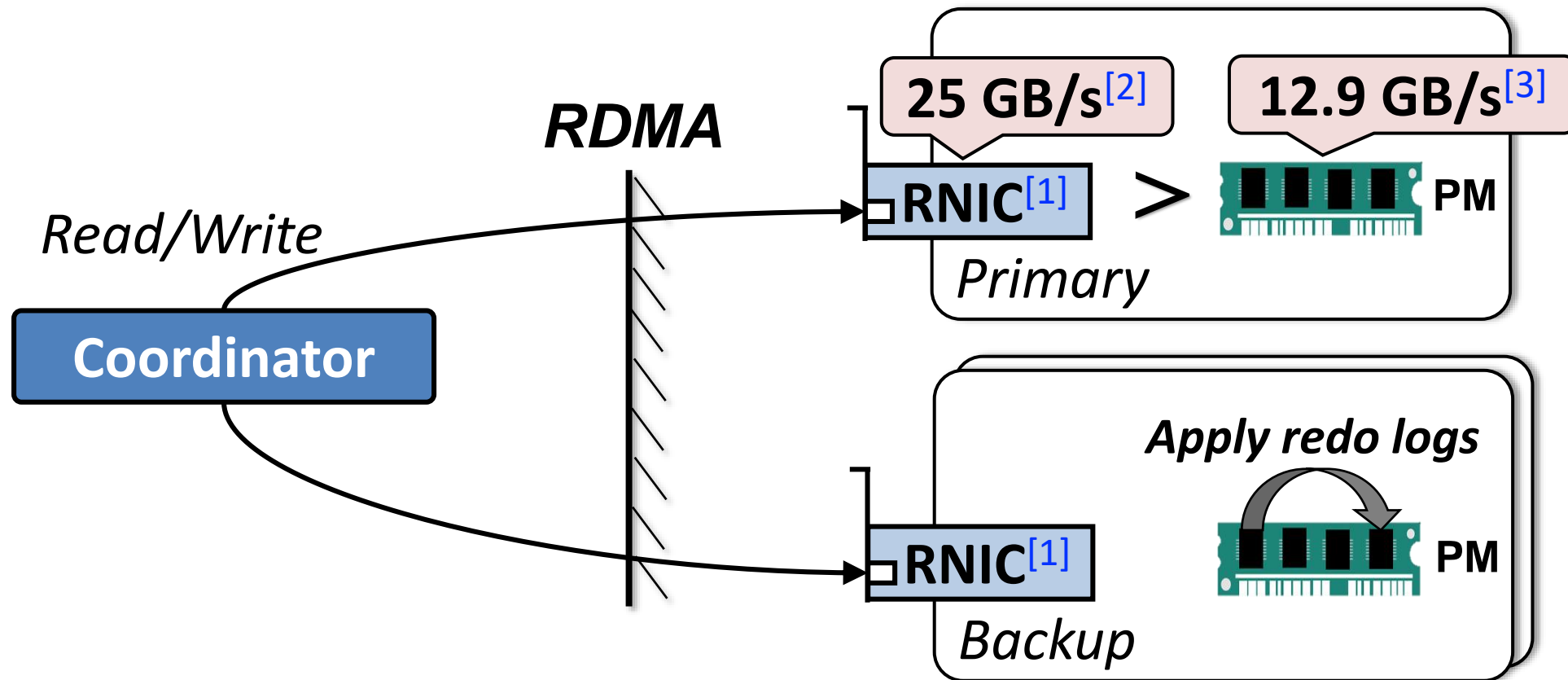


¹ RDMA NIC ² Mellanox ConnectX-6/dual-port ConnectX-5 VPI Card

³ Six interleaved 256GB Intel Optane PM [Jian Yang, et al.@FAST'20]

Challenge 2

- Limited PM bandwidth: *High loads on primary*

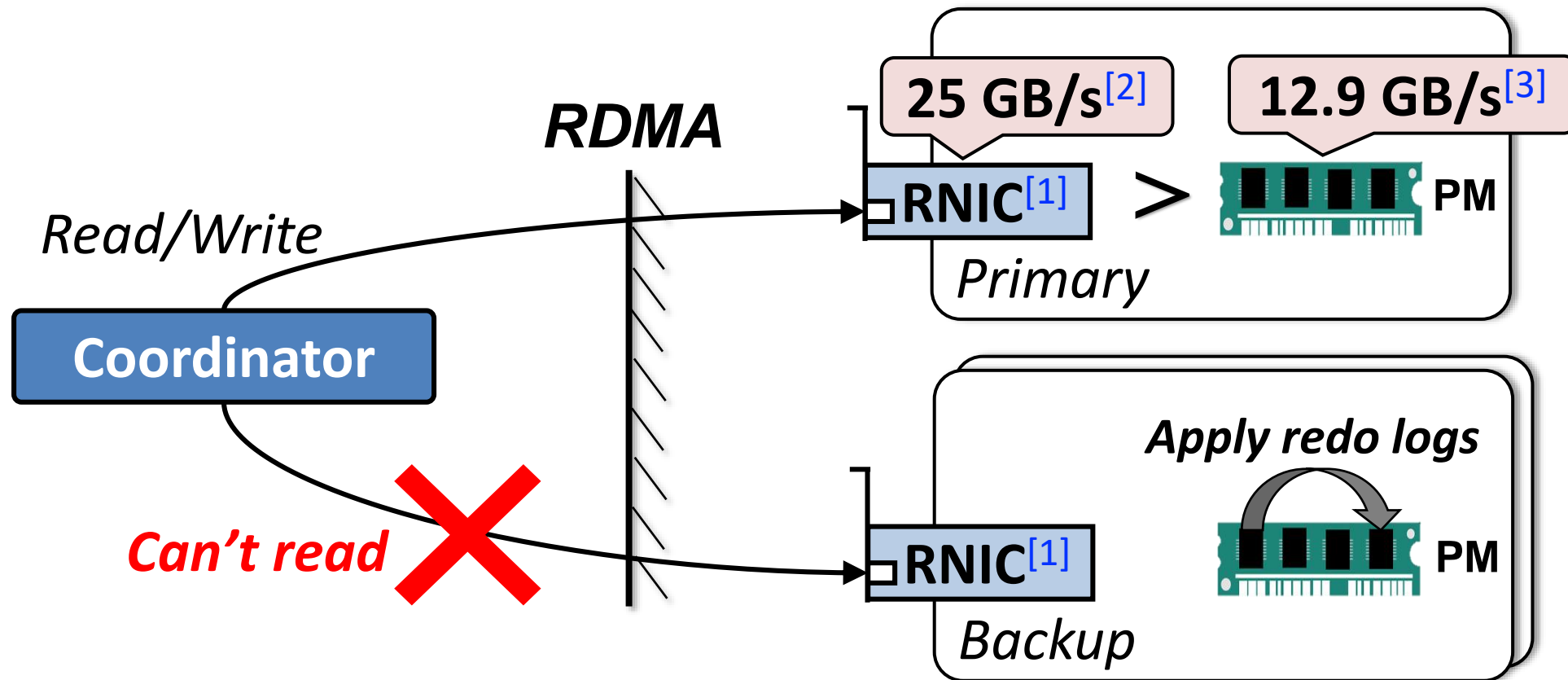


¹ RDMA NIC ² Mellanox ConnectX-6/dual-port ConnectX-5 VPI Card

³ Six interleaved 256GB Intel Optane PM [Jian Yang, et al.@FAST'20]

Challenge 2

- Limited PM bandwidth: *High loads on primary*

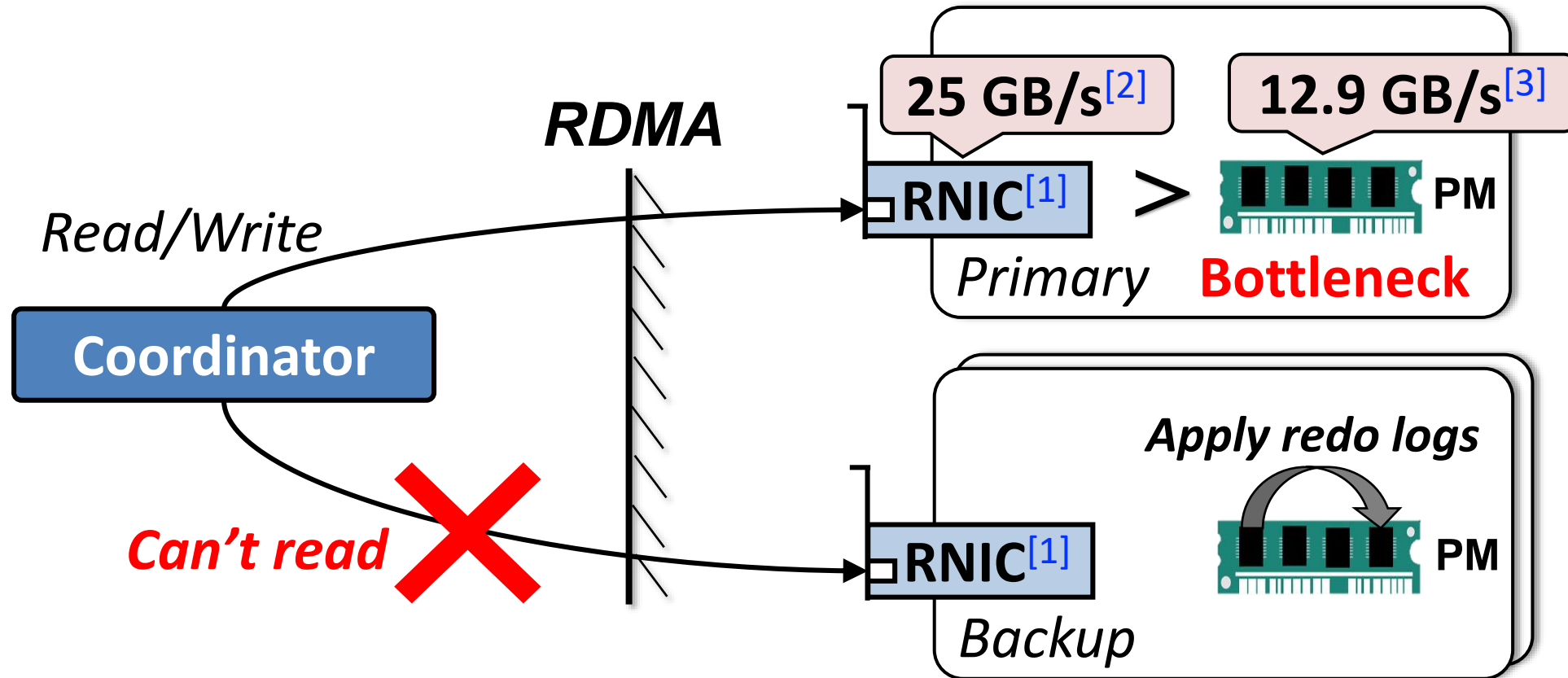


¹ RDMA NIC ² Mellanox ConnectX-6/dual-port ConnectX-5 VPI Card

³ Six interleaved 256GB Intel Optane PM [Jian Yang, et al. @FAST'20]

Challenge 2

- Limited PM bandwidth: *High loads on primary*

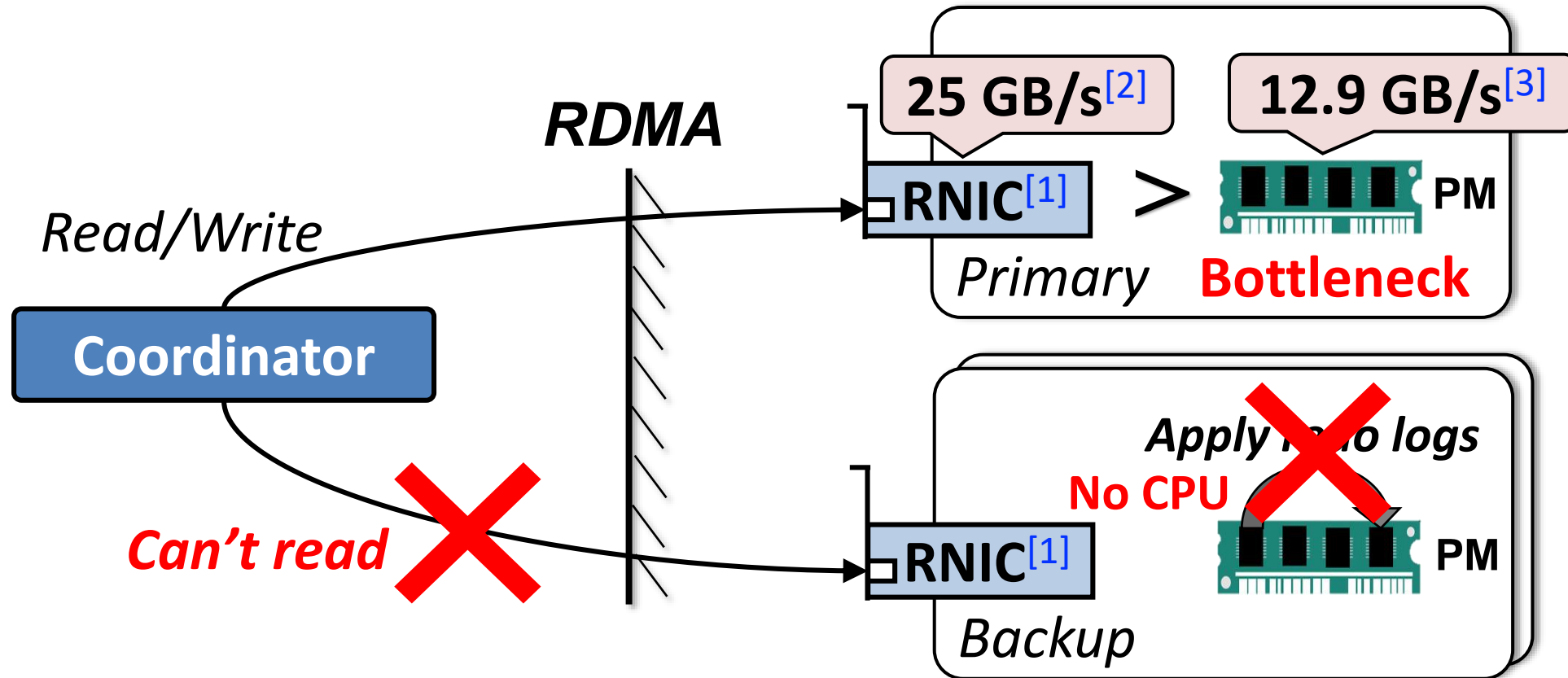


¹ RDMA NIC ² Mellanox ConnectX-6/dual-port ConnectX-5 VPI Card

³ Six interleaved 256GB Intel Optane PM [Jian Yang, et al.@FAST'20]

Challenge 2

- Limited PM bandwidth: *High loads on primary*

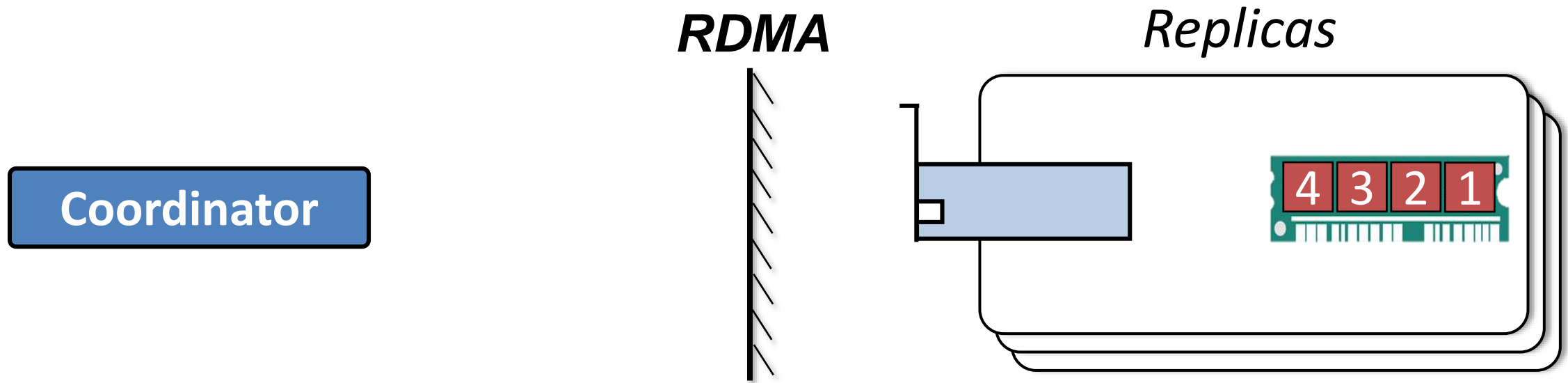


¹ RDMA NIC ² Mellanox ConnectX-6/dual-port ConnectX-5 VPI Card

³ Six interleaved 256GB Intel Optane PM [Jian Yang, et al. @FAST'20]

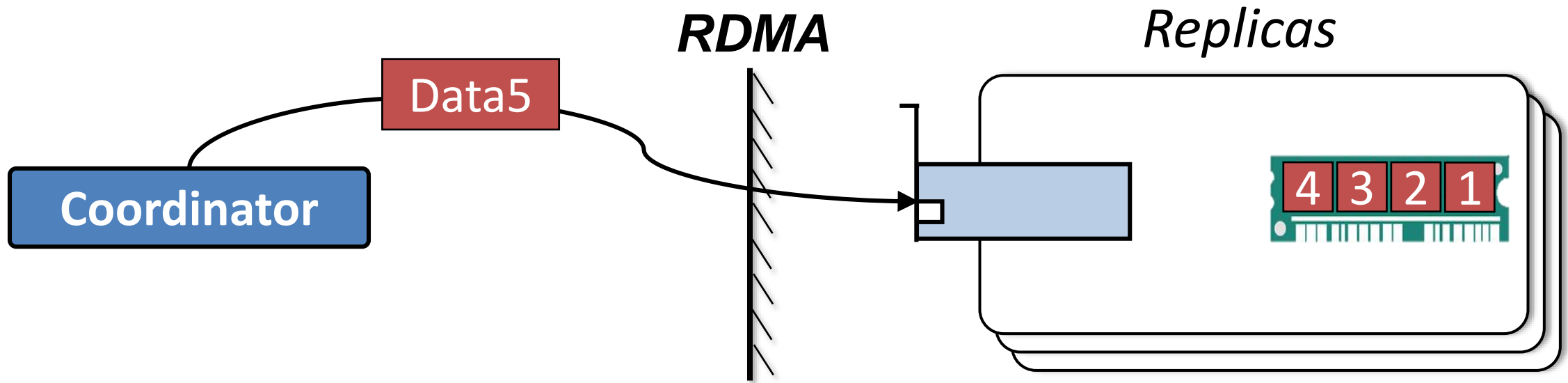
Challenge 3

- Lack of remote persistency guarantee: *Inconsistent write*



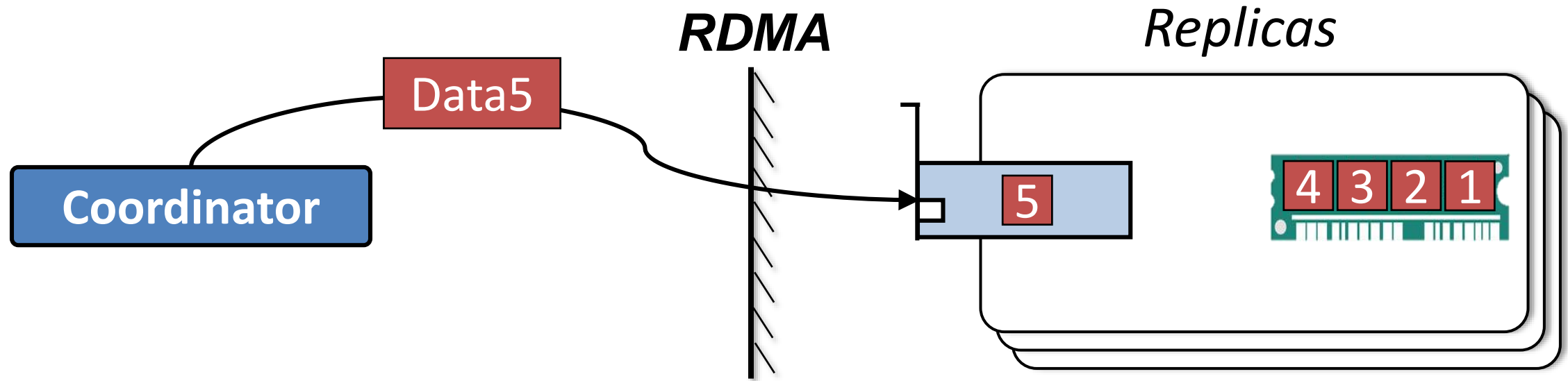
Challenge 3

- Lack of remote persistency guarantee: *Inconsistent write*



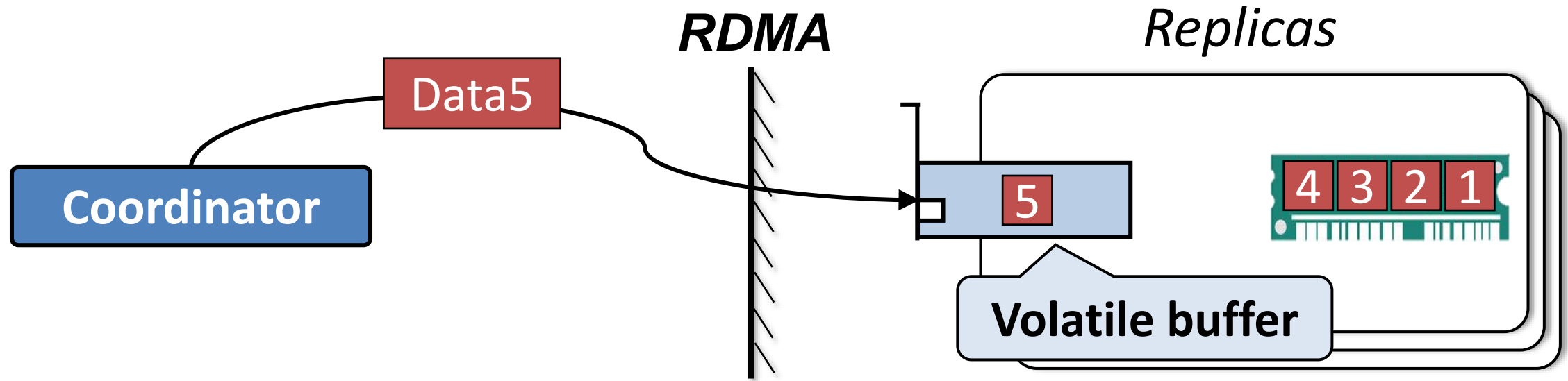
Challenge 3

- Lack of remote persistency guarantee: *Inconsistent write*



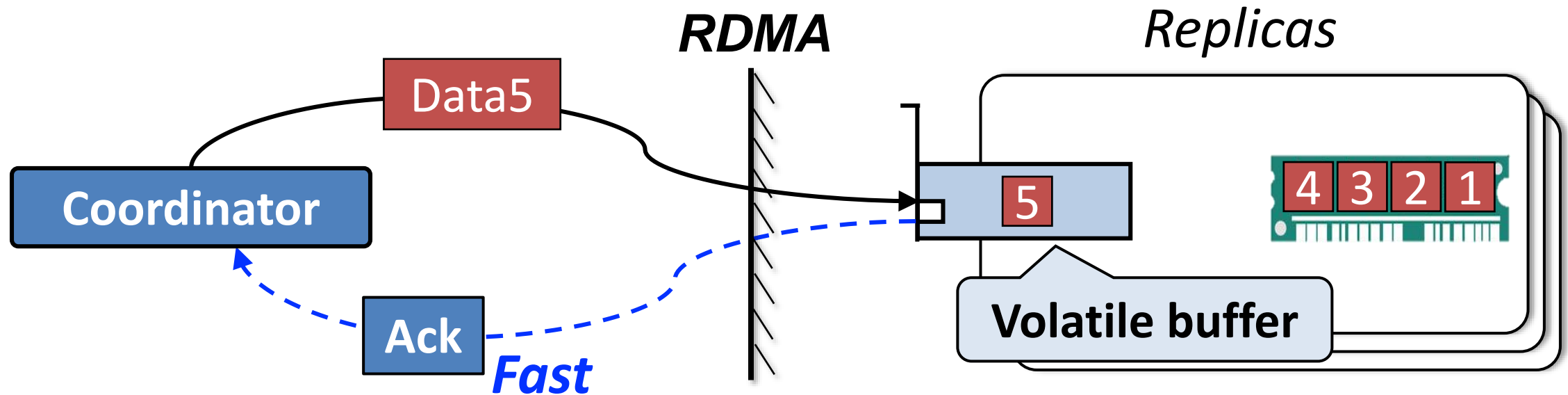
Challenge 3

- Lack of remote persistency guarantee: *Inconsistent write*



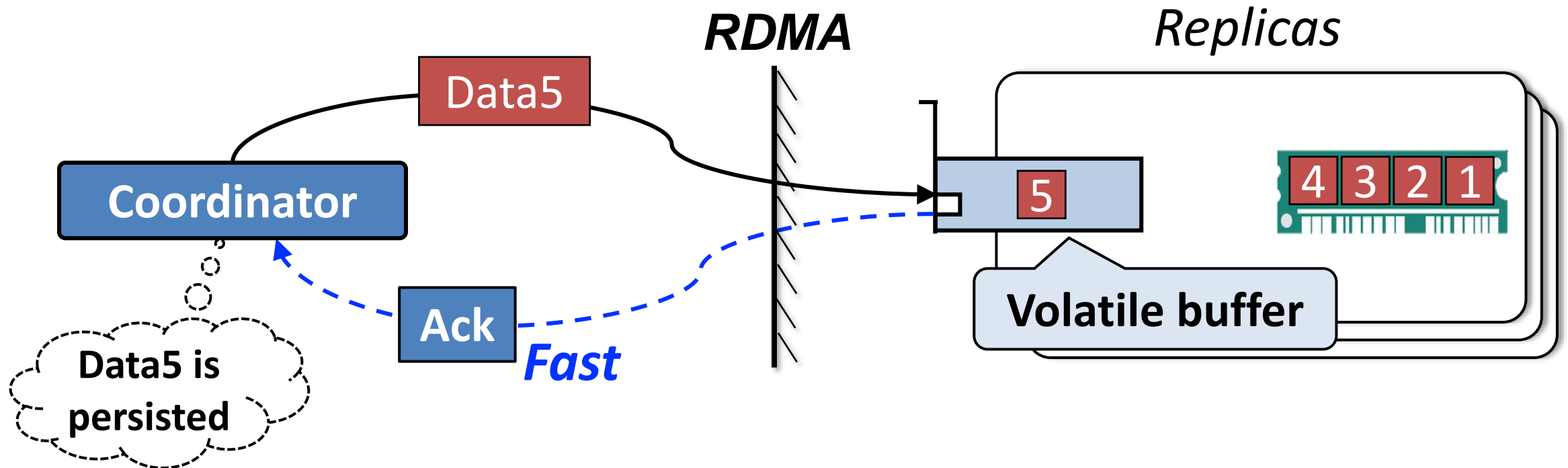
Challenge 3

- Lack of remote persistency guarantee: *Inconsistent write*



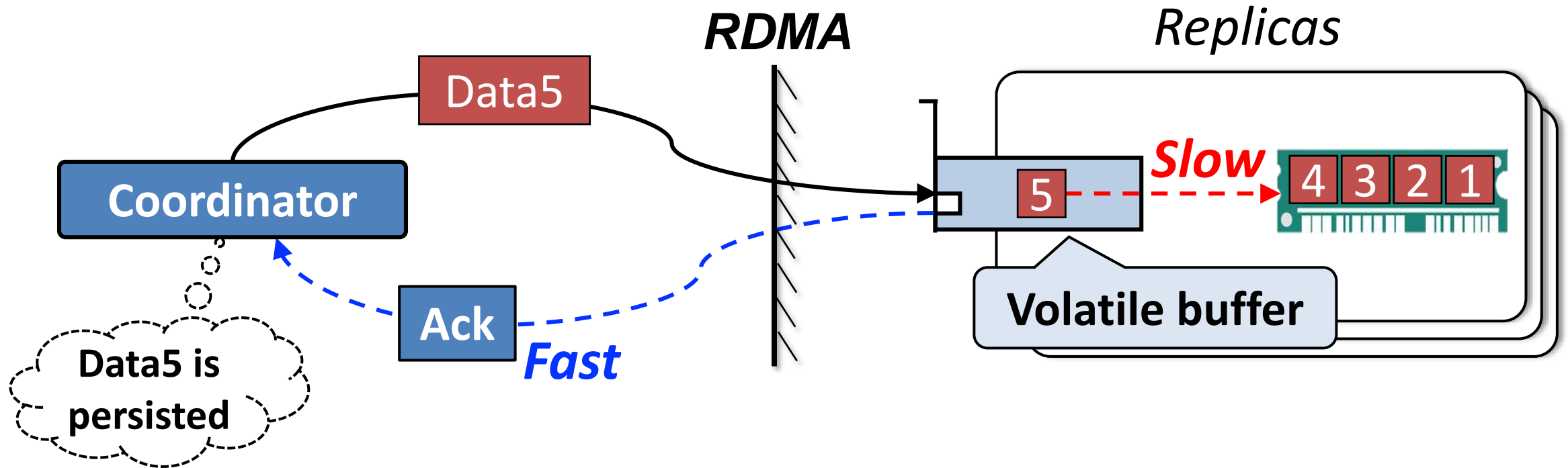
Challenge 3

- Lack of remote persistency guarantee: *Inconsistent write*



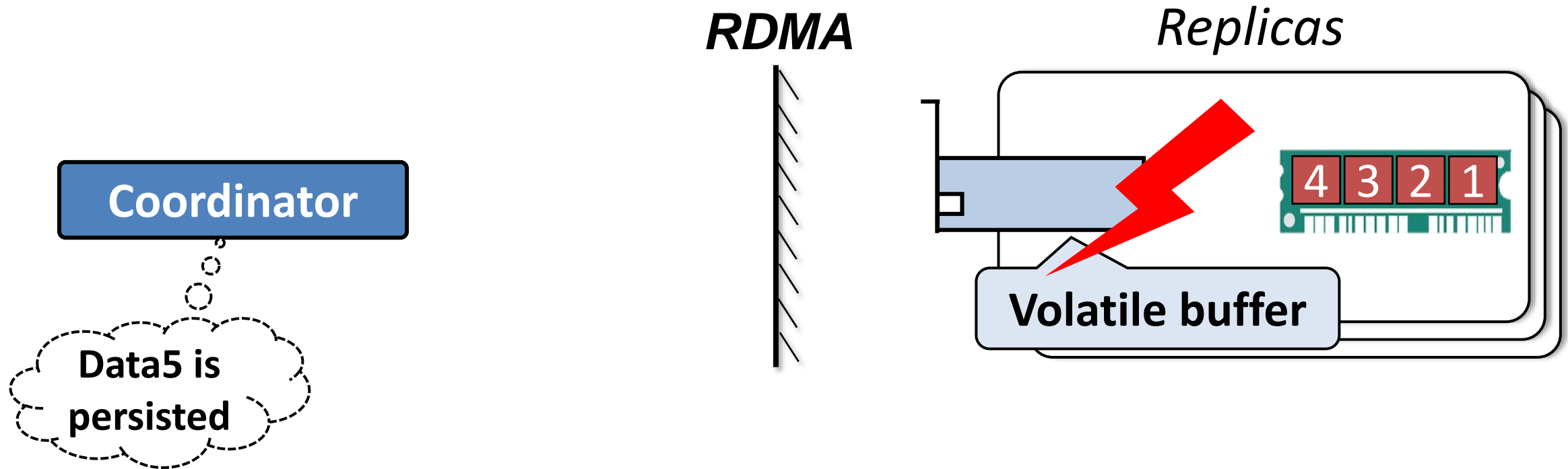
Challenge 3

- Lack of remote persistency guarantee: *Inconsistent write*



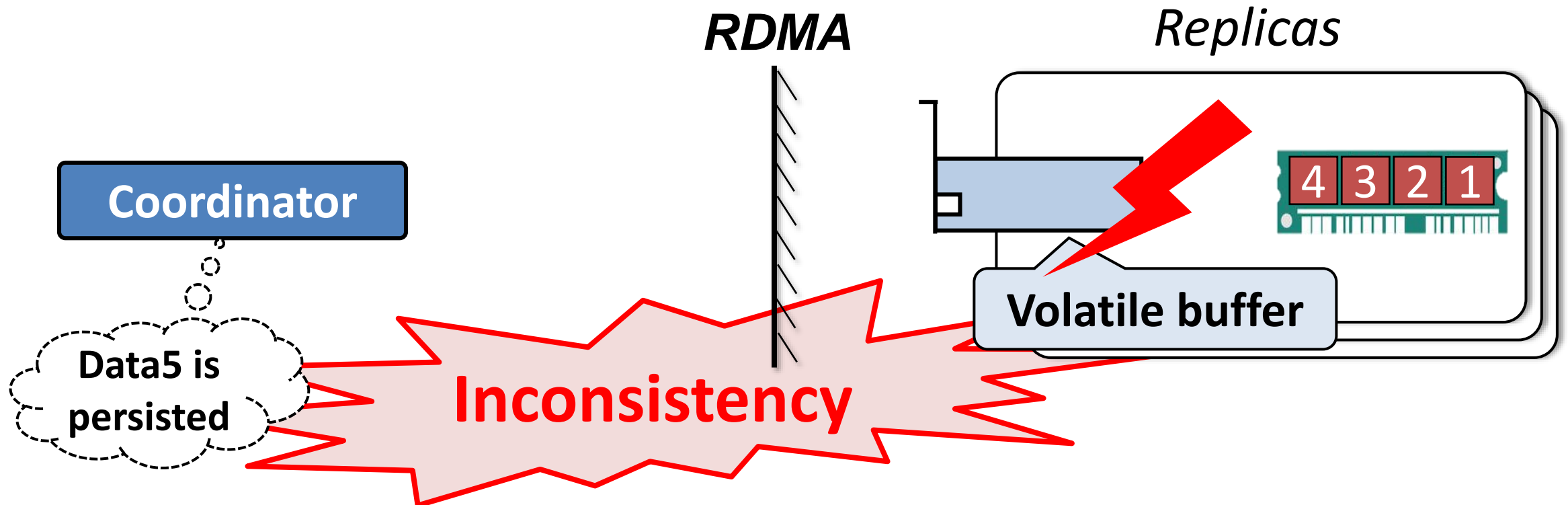
Challenge 3

- Lack of remote persistency guarantee: *Inconsistent write*



Challenge 3

- Lack of remote persistency guarantee: *Inconsistent write*



FORD: Solutions

- Long-latency processing: *Many round trips*
- Limited PM bandwidth: *High loads on primary*
- Lack of remote persistency guarantee: *Inconsistent write*

FORD: Solutions

➤ Long-latency processing: *Many round trips*

- Hitchhiked Locking
- Coalescent Commit

➤ Reduce round trips to decrease latency

➤ Limited PM bandwidth: *High loads on primary*

➤ Lack of remote persistency guarantee: *Inconsistent write*

FORD: Solutions

➤ Long-latency processing: *Many round trips*

- Hitchhiked Locking
- Coalescent Commit

➡ Reduce round trips to decrease latency

➤ Limited PM bandwidth: *High loads on primary*

- Backup-enabled Read

➡ Balance load to improve throughput

➤ Lack of remote persistency guarantee: *Inconsistent write*

FORD: Solutions

➤ Long-latency processing: *Many round trips*

- Hitchhiked Locking
- Coalescent Commit

➡ Reduce round trips to decrease latency

➤ Limited PM bandwidth: *High loads on primary*

- Backup-enabled Read

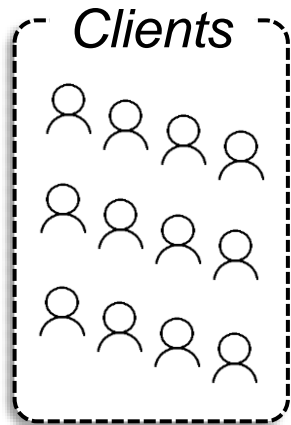
➡ Balance load to improve throughput

➤ Lack of remote persistency guarantee: *Inconsistent write*

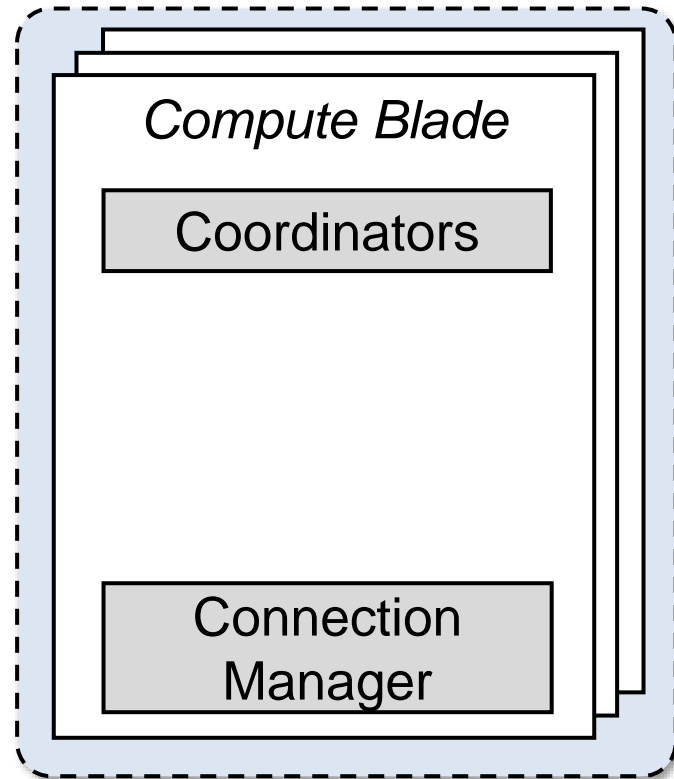
- Selective Remote Flush

➡ Guarantee remote persistency with low overhead

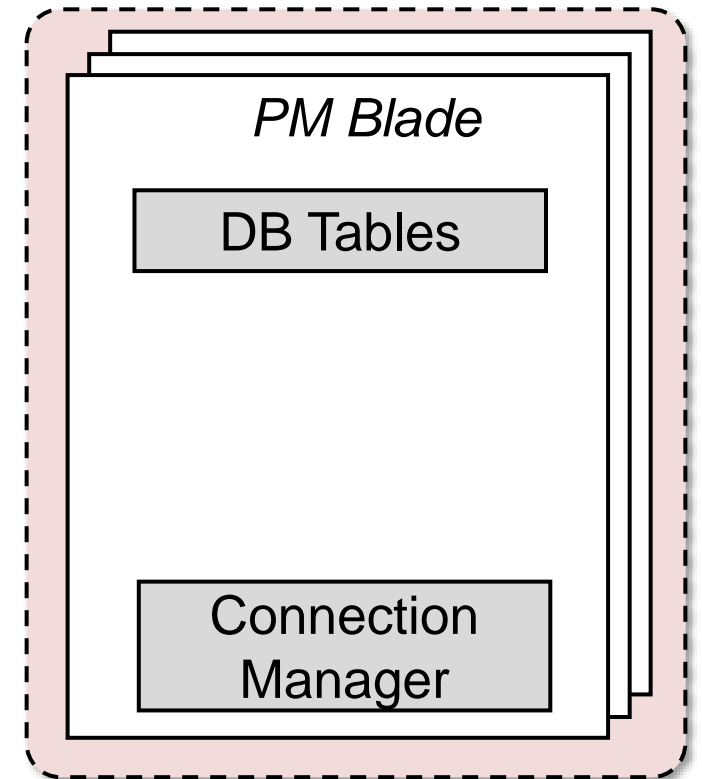
Overview



- Init: **1-2**
- Run: **3-5**



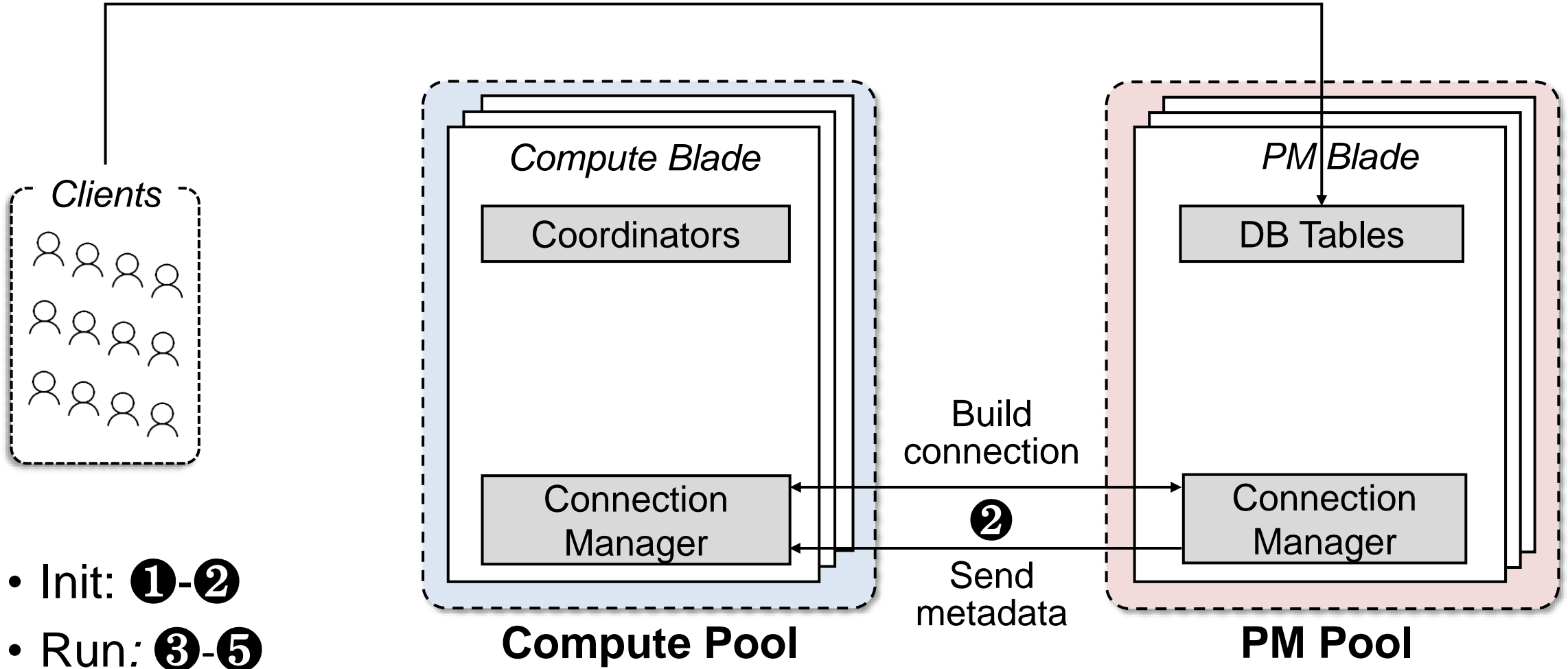
Compute Pool



PM Pool

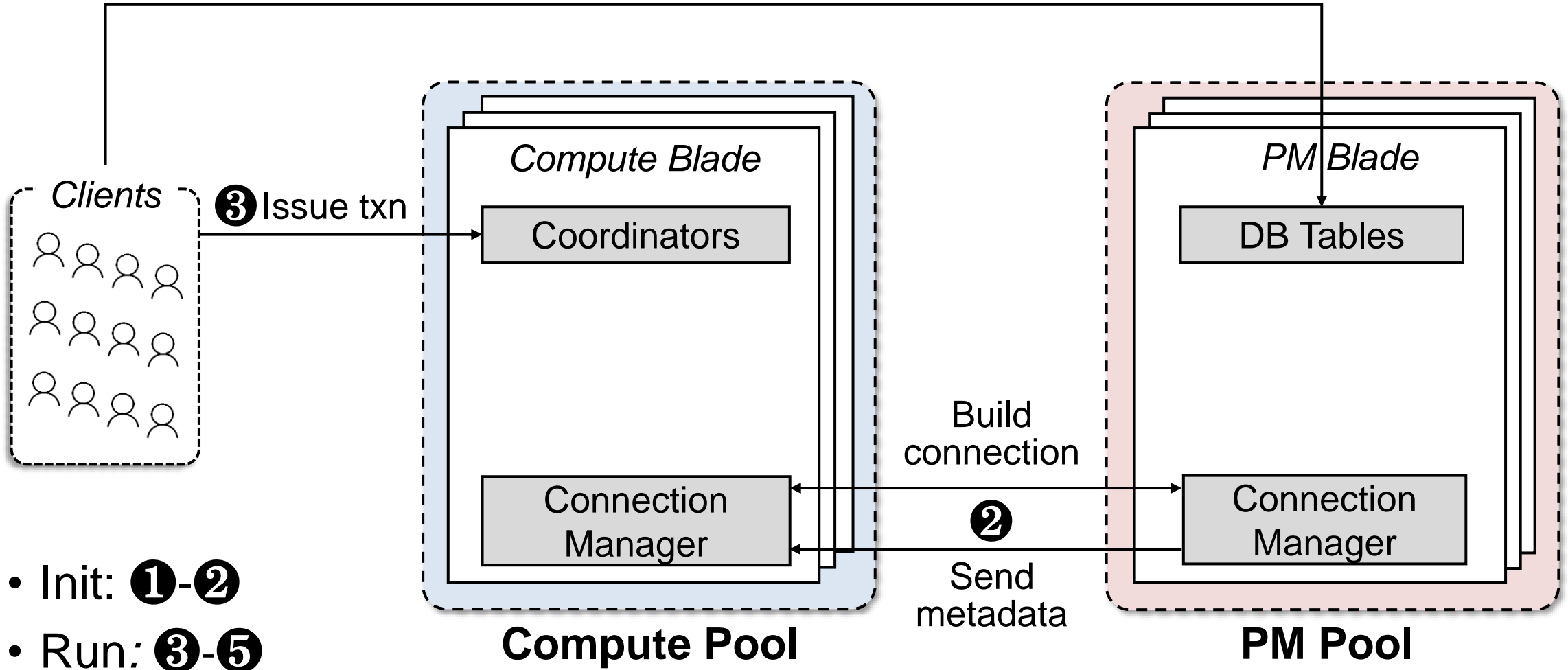
Overview

1 Allocate memory and load table

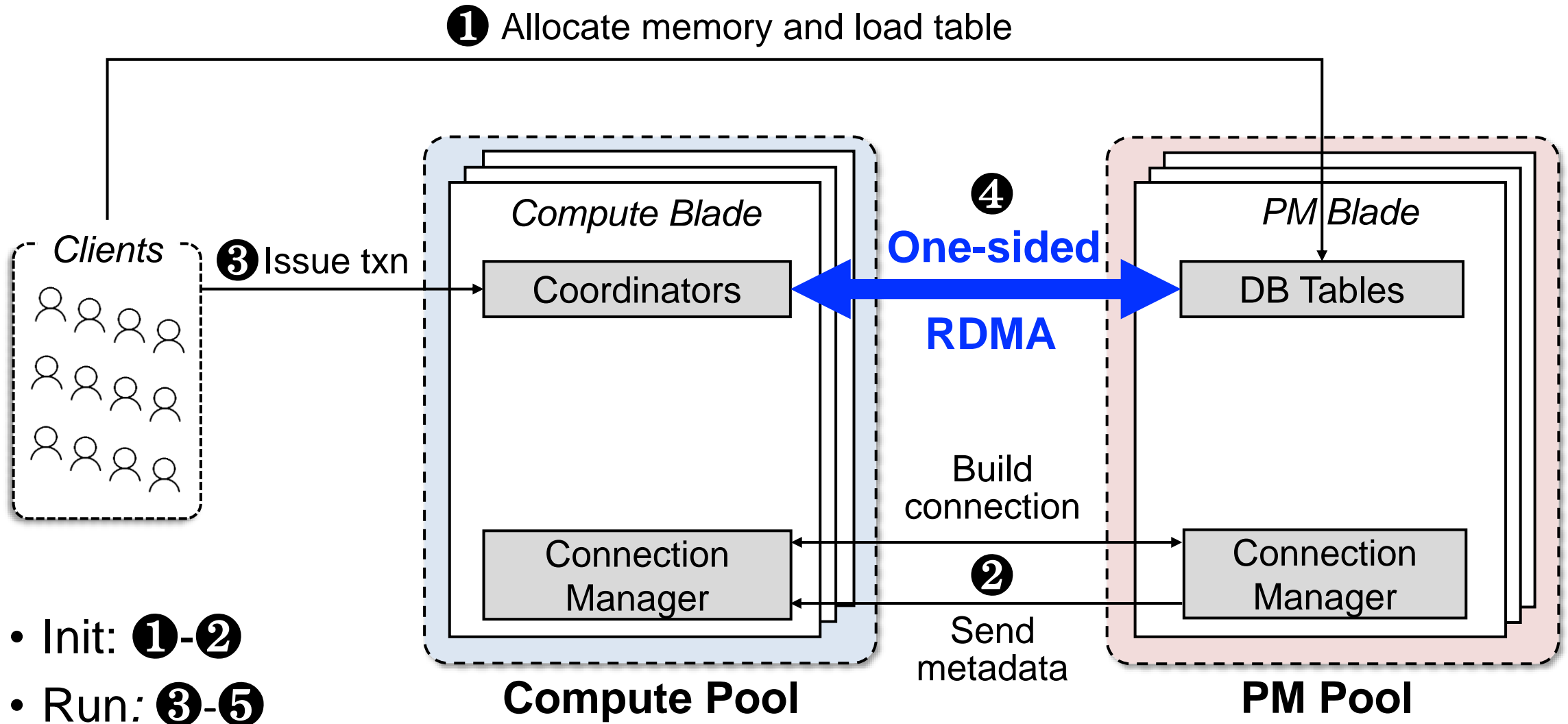


Overview

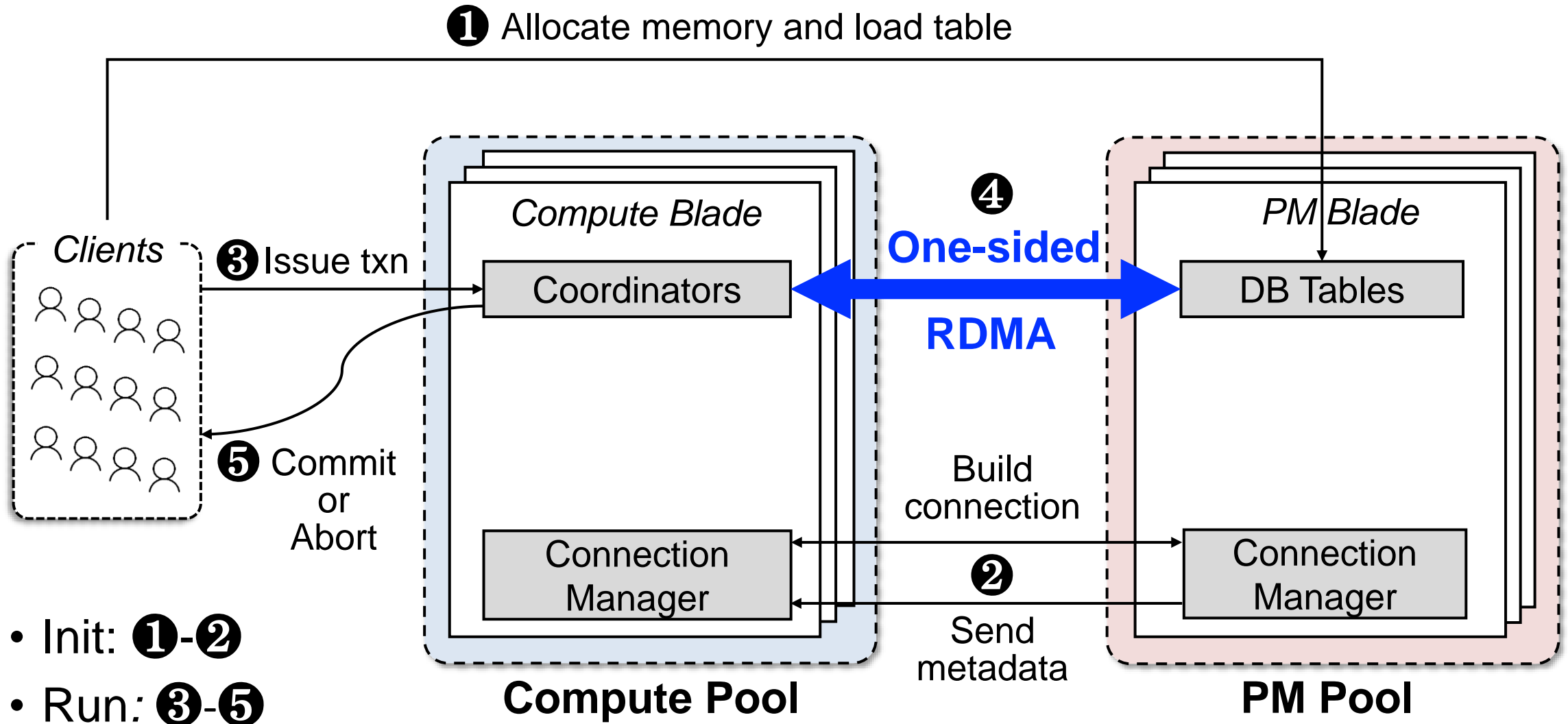
1 Allocate memory and load table



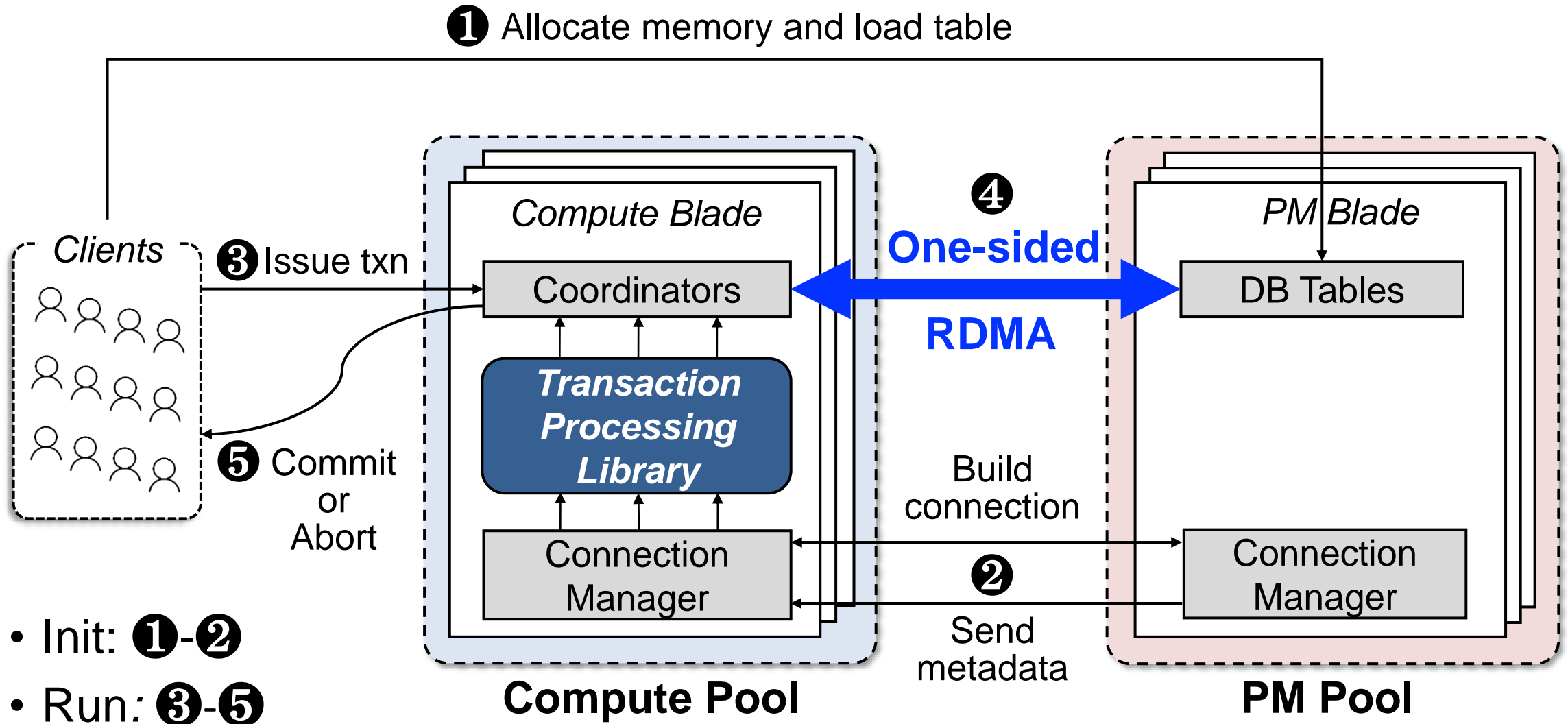
Overview



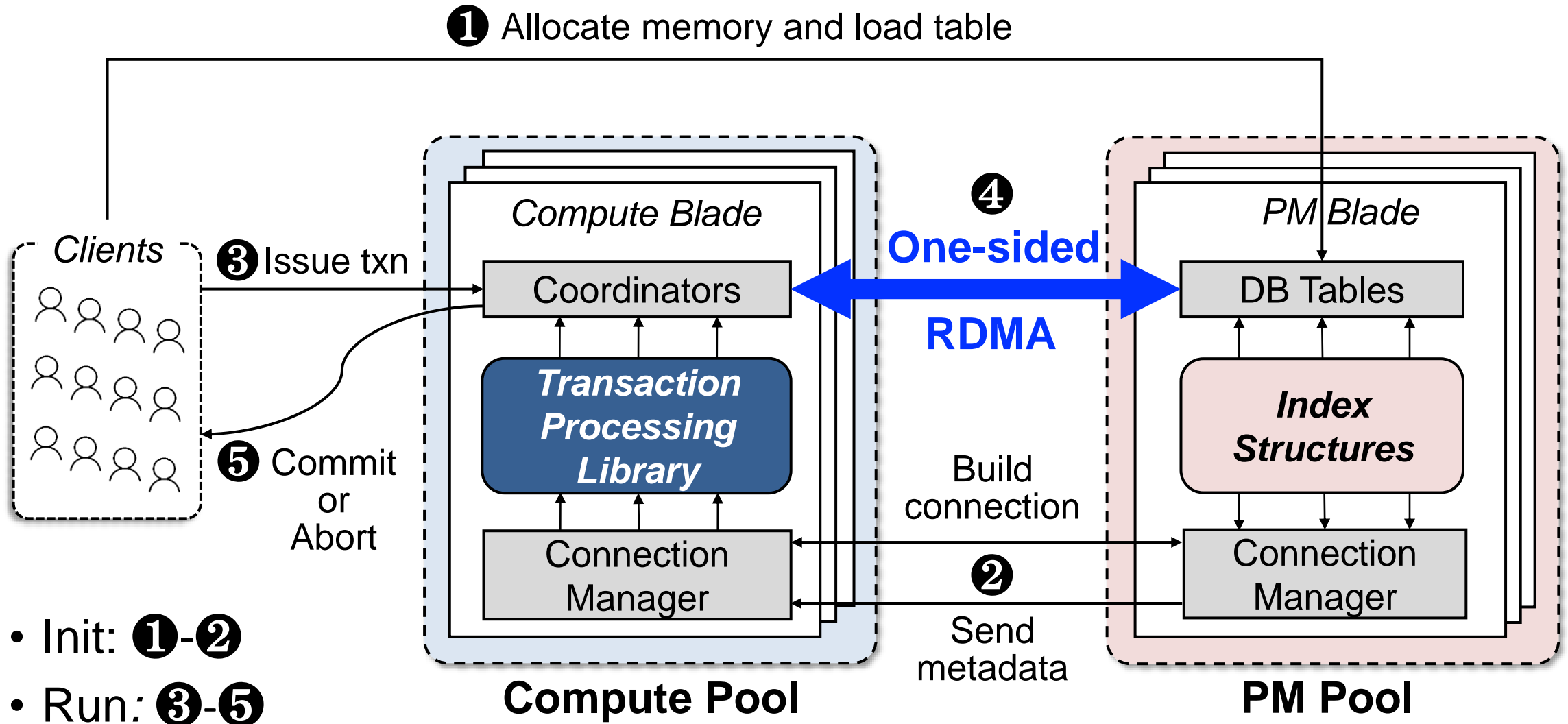
Overview



Overview

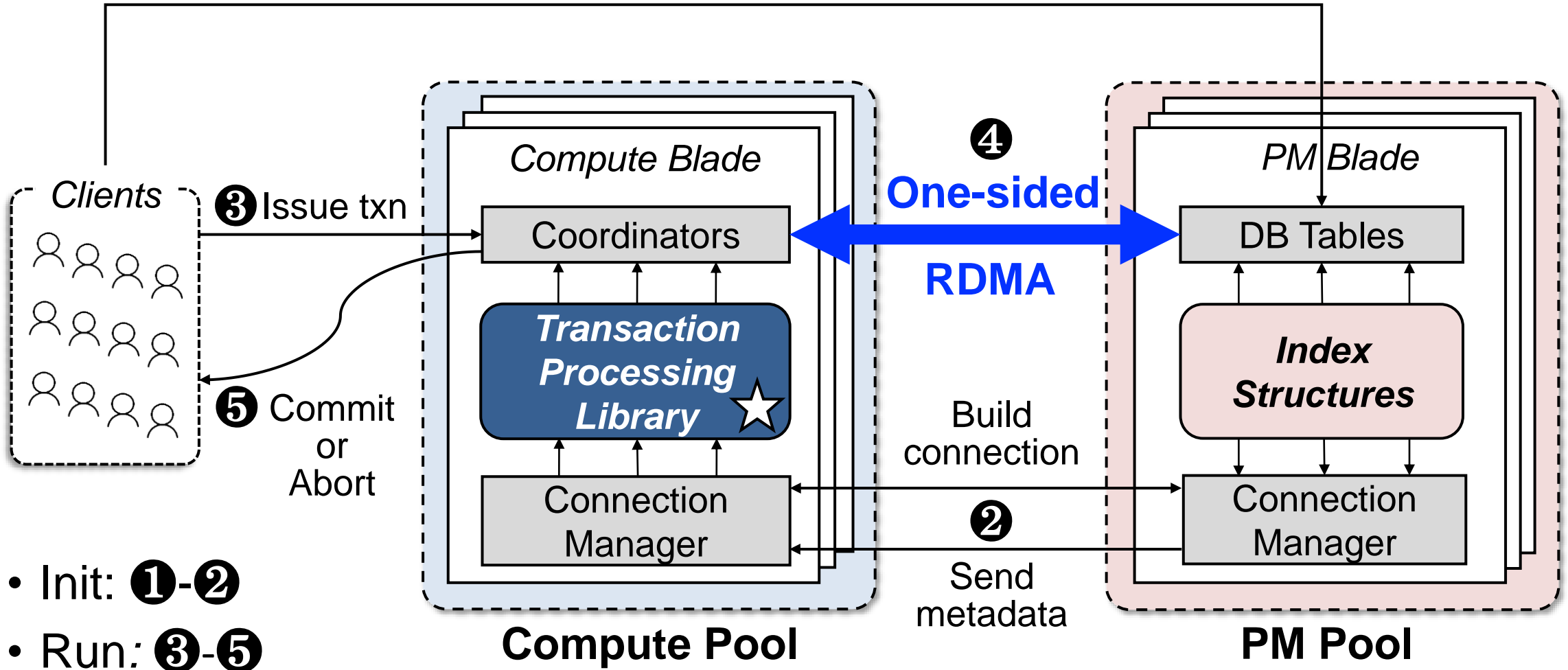


Overview



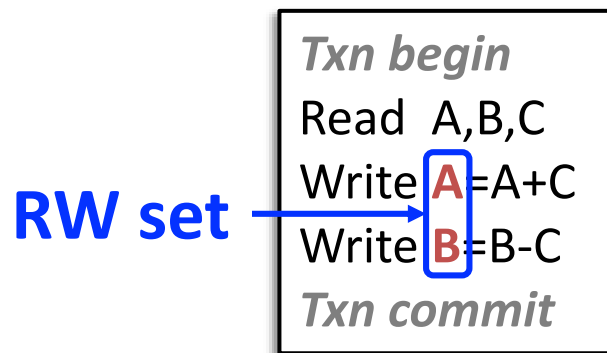
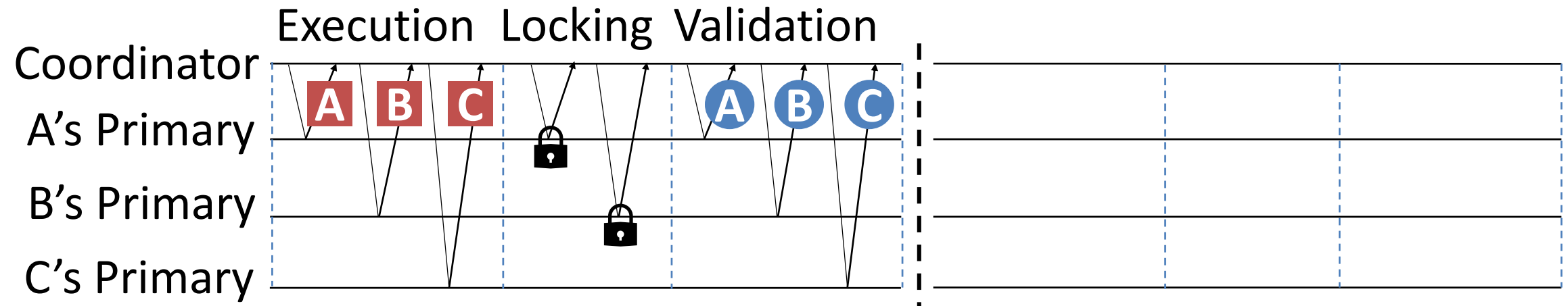
Overview

① Allocate memory and load table



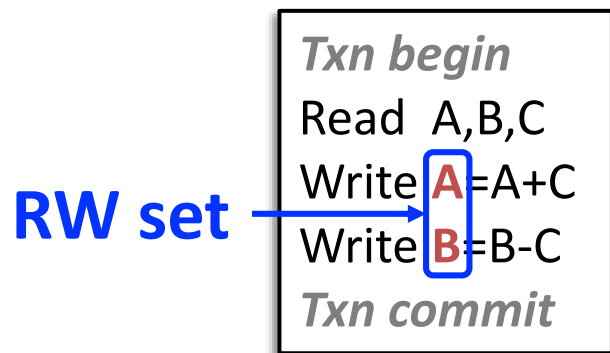
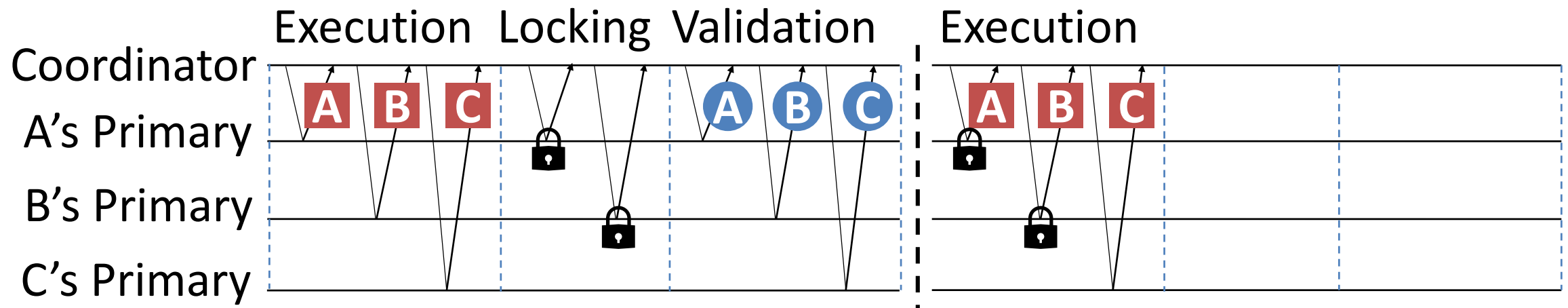
Hitchhiked Locking

- Read and lock the read-write (RW) set in execution
 - Avoid subsequent locking and validations
 - No lock on the read-only data



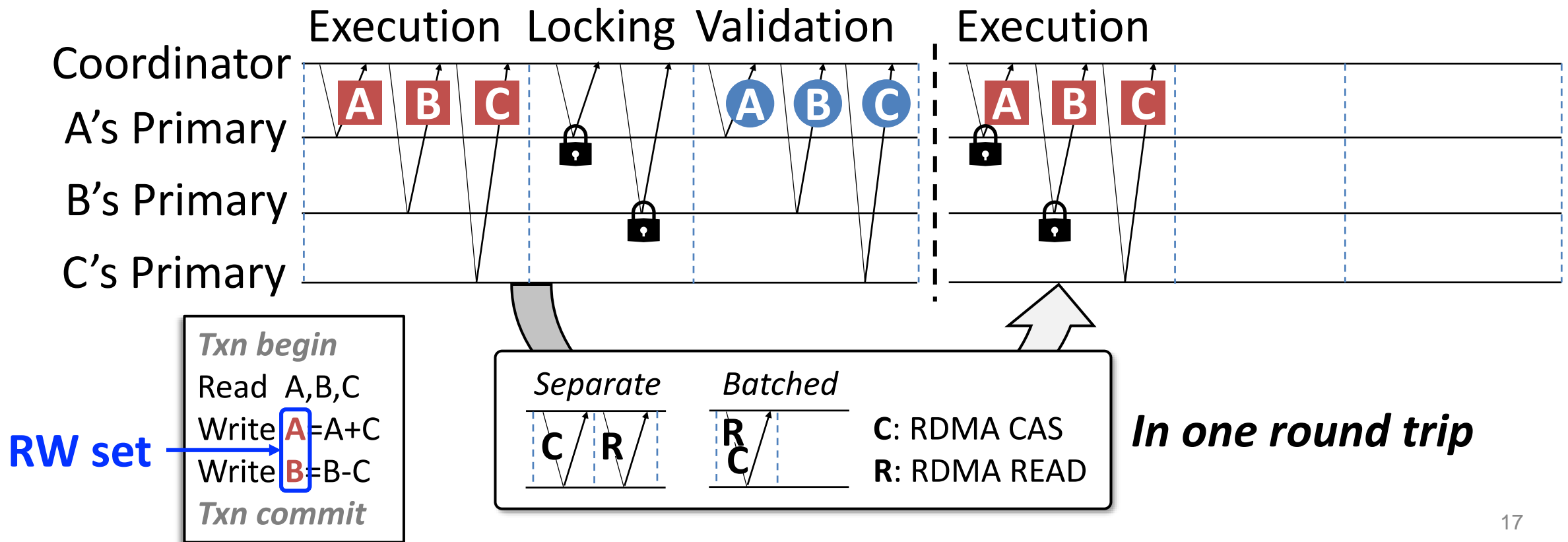
Hitchhiked Locking

- Read and lock the read-write (RW) set in execution
 - Avoid subsequent locking and validations
 - No lock on the read-only data



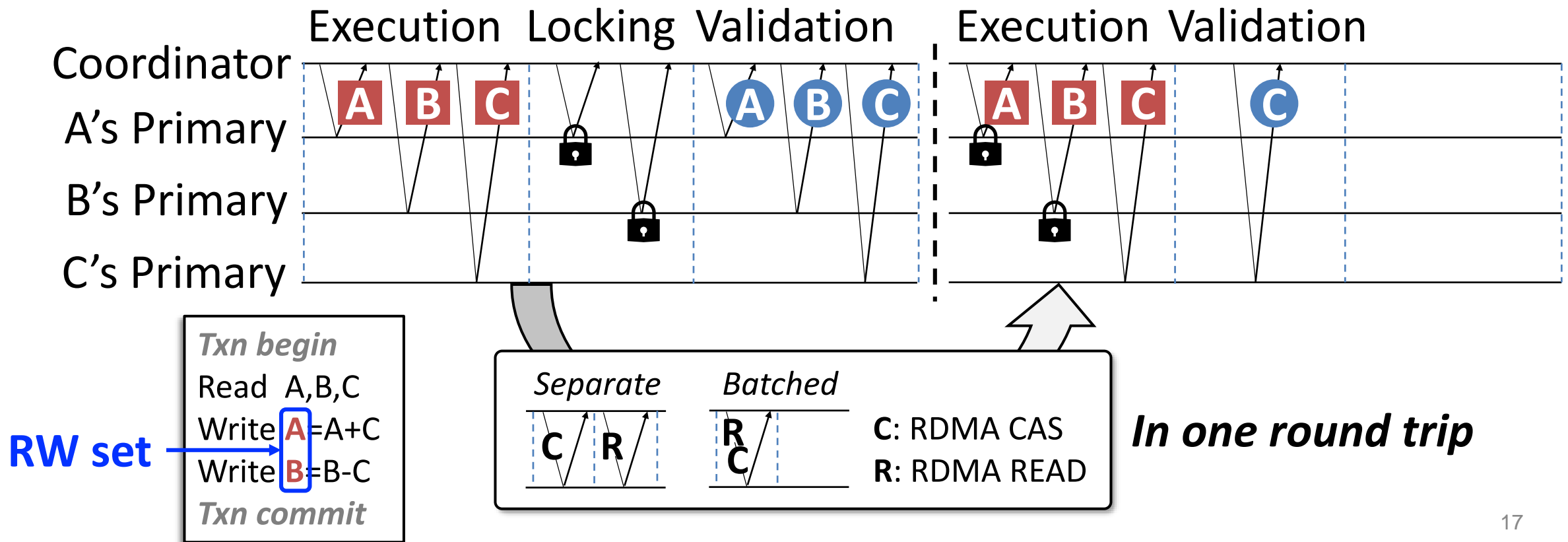
Hitchhiked Locking

- Read and lock the read-write (RW) set in execution
 - Avoid subsequent locking and validations
 - No lock on the read-only data



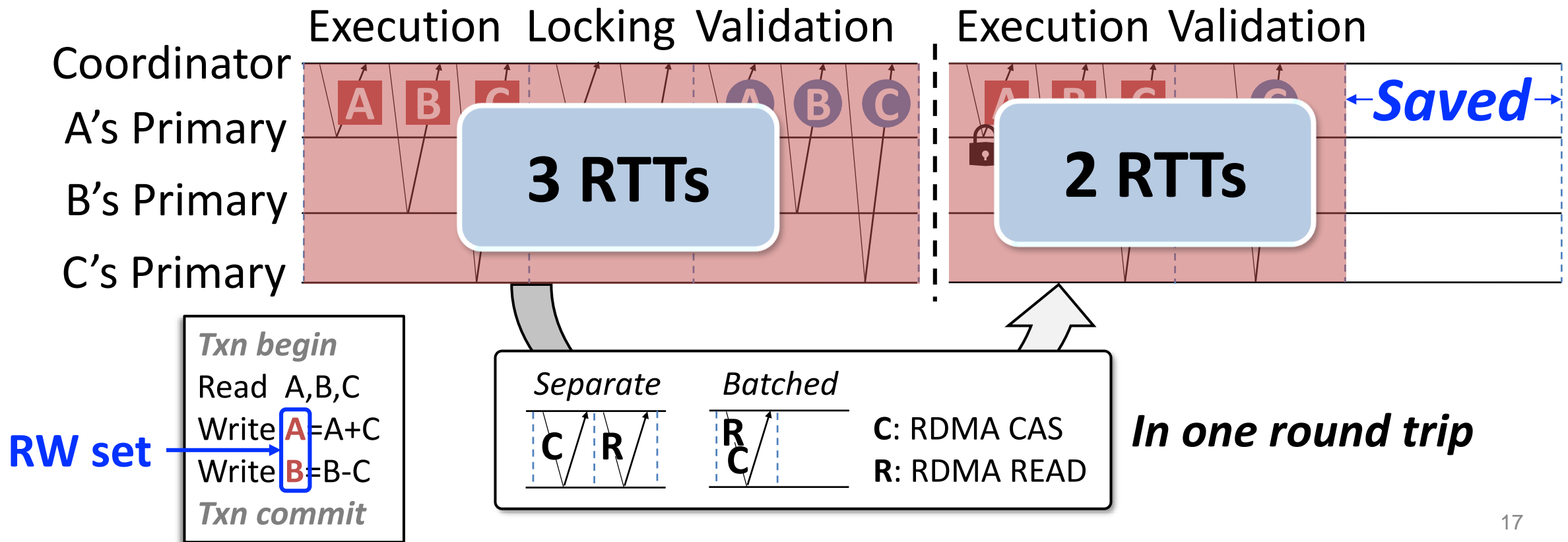
Hitchhiked Locking

- Read and lock the read-write (RW) set in execution
 - Avoid subsequent locking and validations
 - No lock on the read-only data



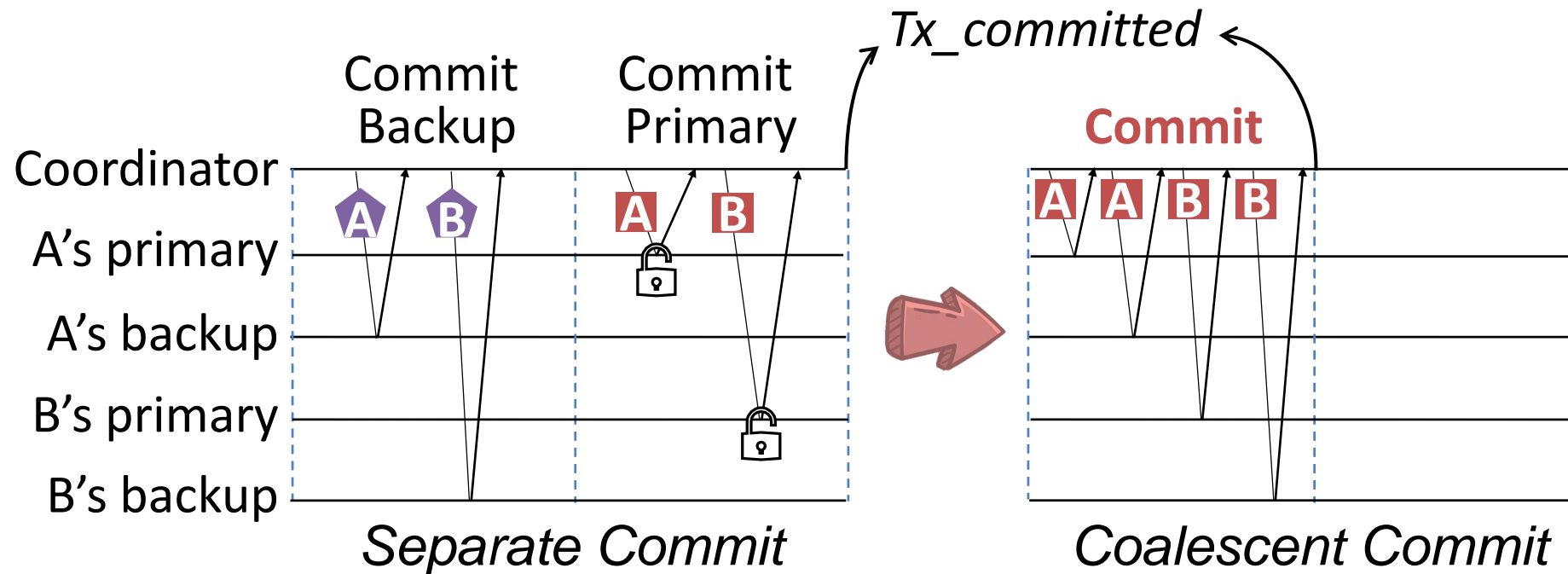
Hitchhiked Locking

- Read and lock the read-write (RW) set in execution
 - Avoid subsequent locking and validations
 - No lock on the read-only data



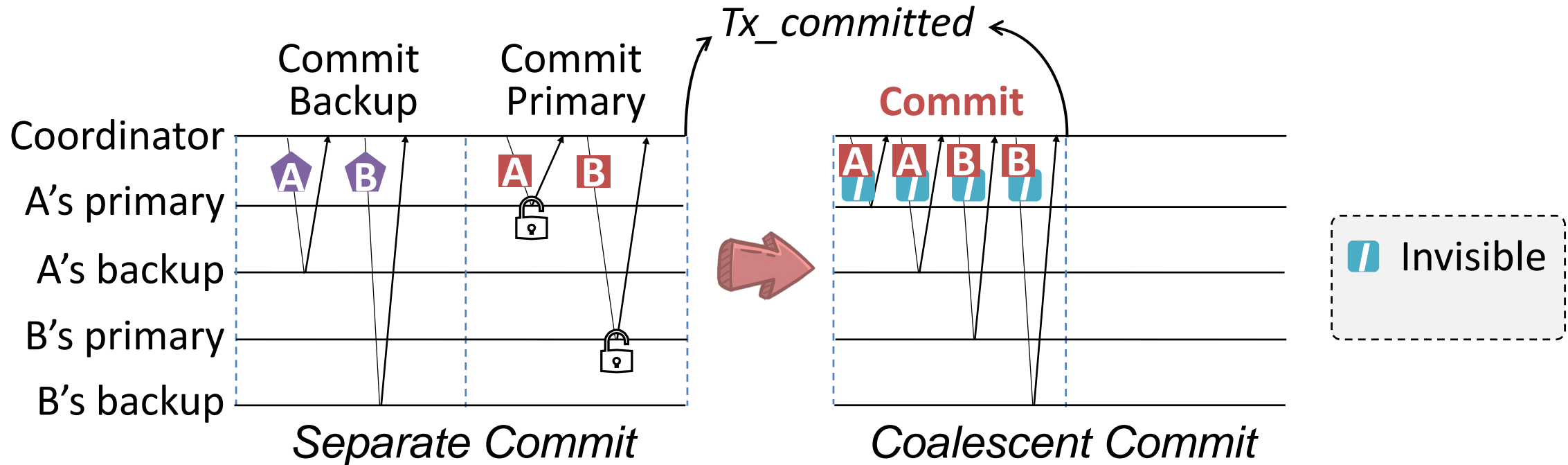
Coalescent Commit

- Commit all replicas **together** in one round trip
 - In-place update: Parallel undo logging in execution phase
 - Prevent reading partial updates: Control data visibility



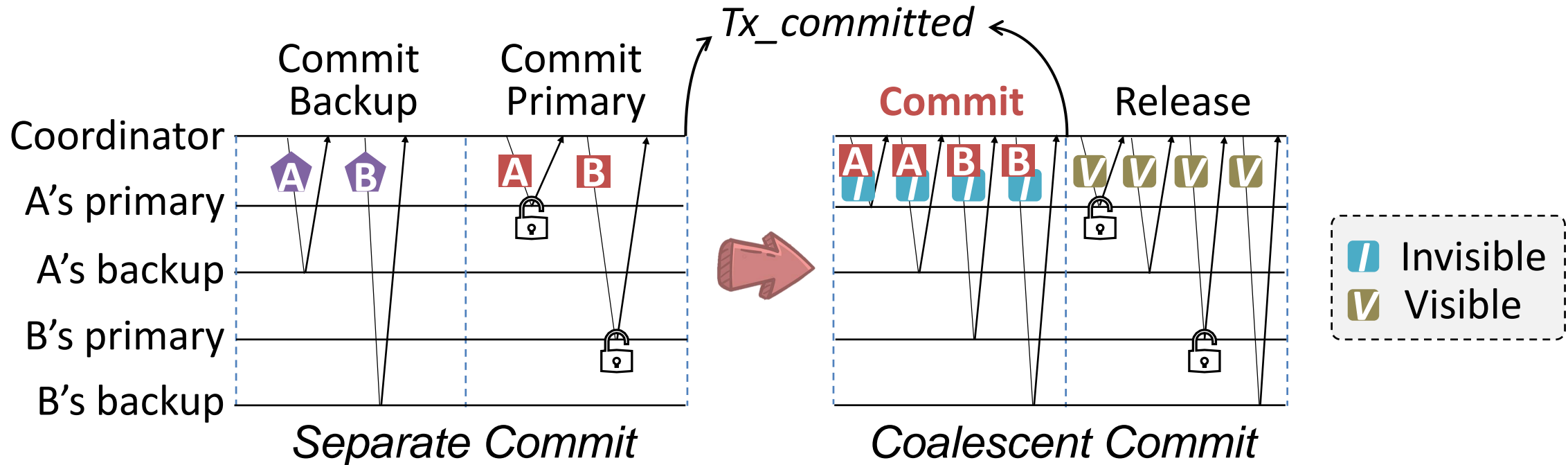
Coalescent Commit

- Commit all replicas **together** in one round trip
 - In-place update: Parallel undo logging in execution phase
 - Prevent reading partial updates: Control data visibility



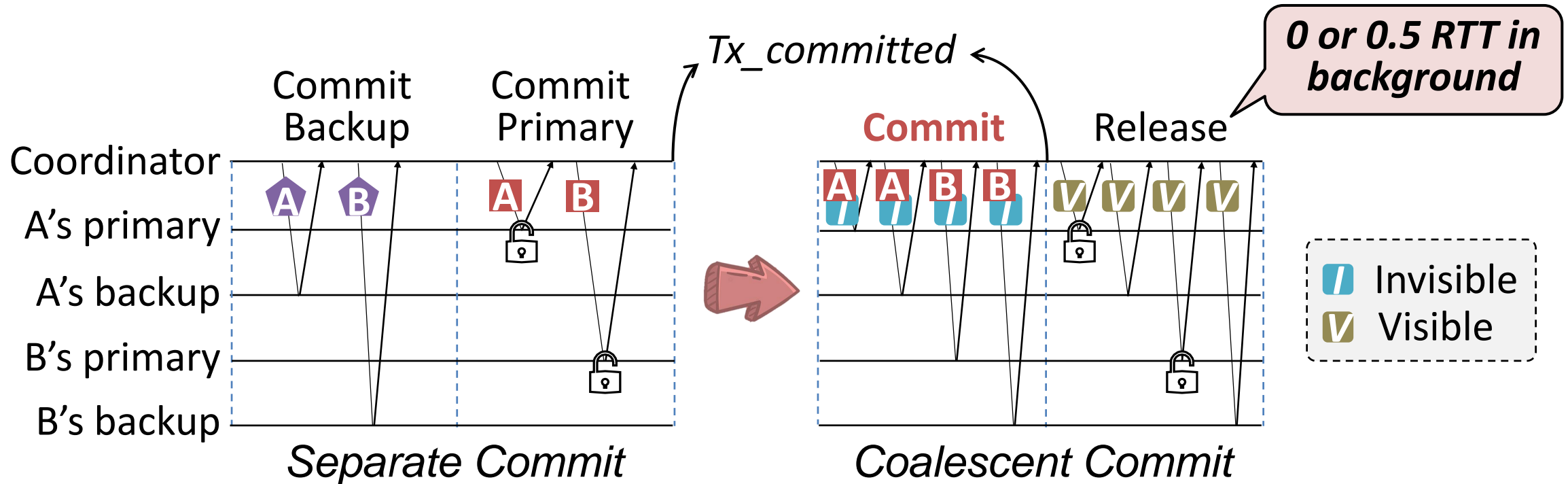
Coalescent Commit

- Commit all replicas **together** in one round trip
 - In-place update: Parallel undo logging in execution phase
 - Prevent reading partial updates: Control data visibility



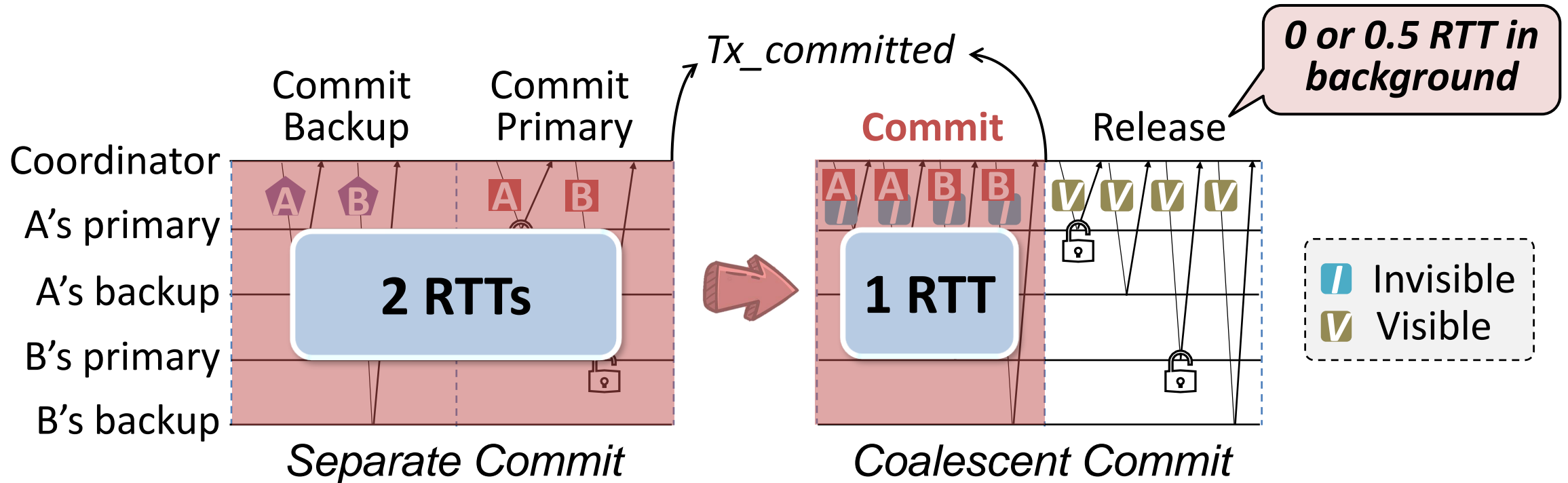
Coalescent Commit

- Commit all replicas **together** in one round trip
 - In-place update: Parallel undo logging in execution phase
 - Prevent reading partial updates: Control data visibility



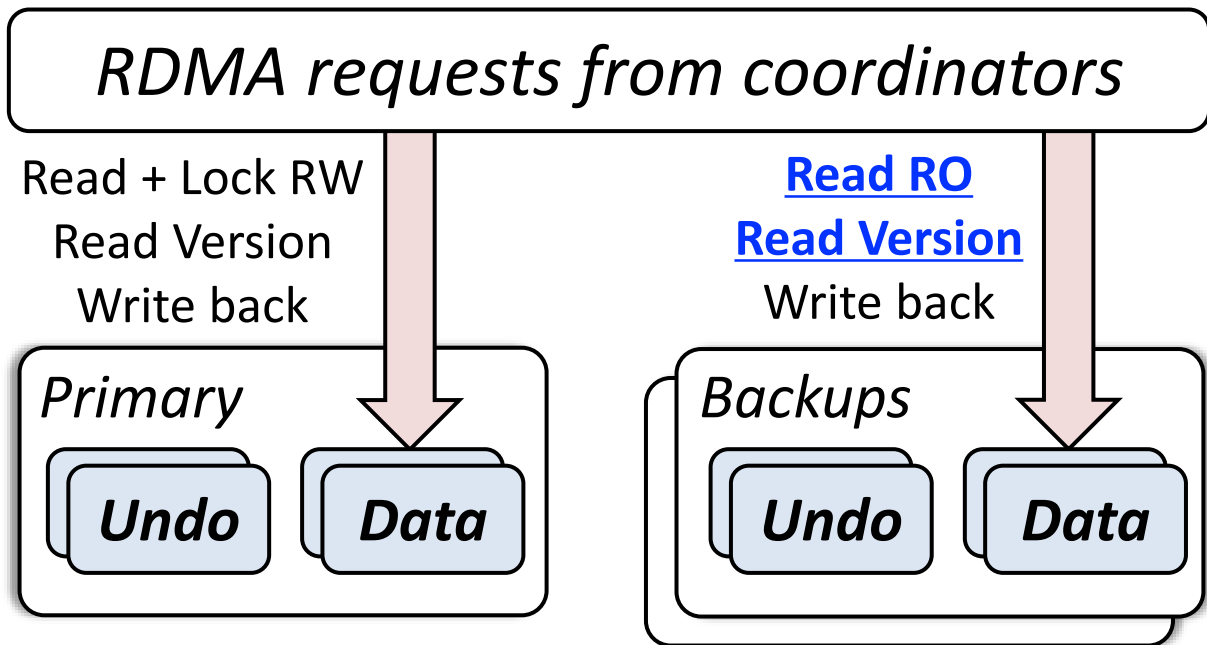
Coalescent Commit

- Commit all replicas **together** in one round trip
 - In-place update: Parallel undo logging in execution phase
 - Prevent reading partial updates: Control data visibility



Backup-enabled Read

- Allows **backups** to serve **read-only (RO)** data
 - In-place update → No address redirection
 - Undo logging → No data migration → No CPU involvement in PM pool

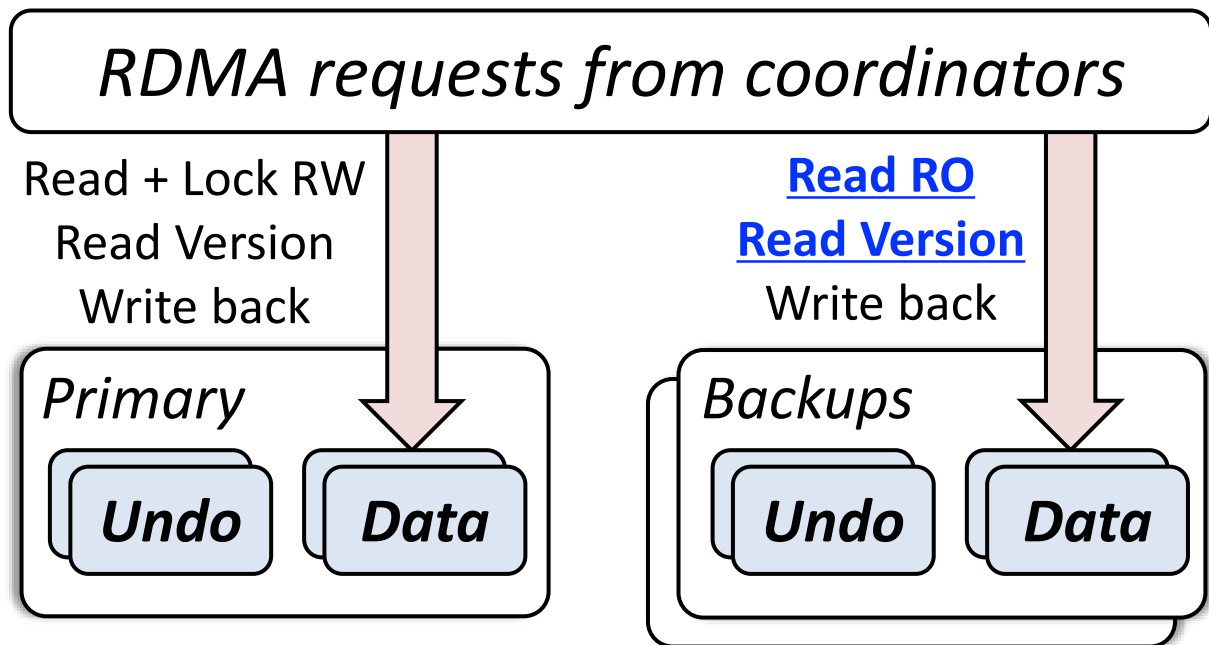


Balance load → Improve throughput

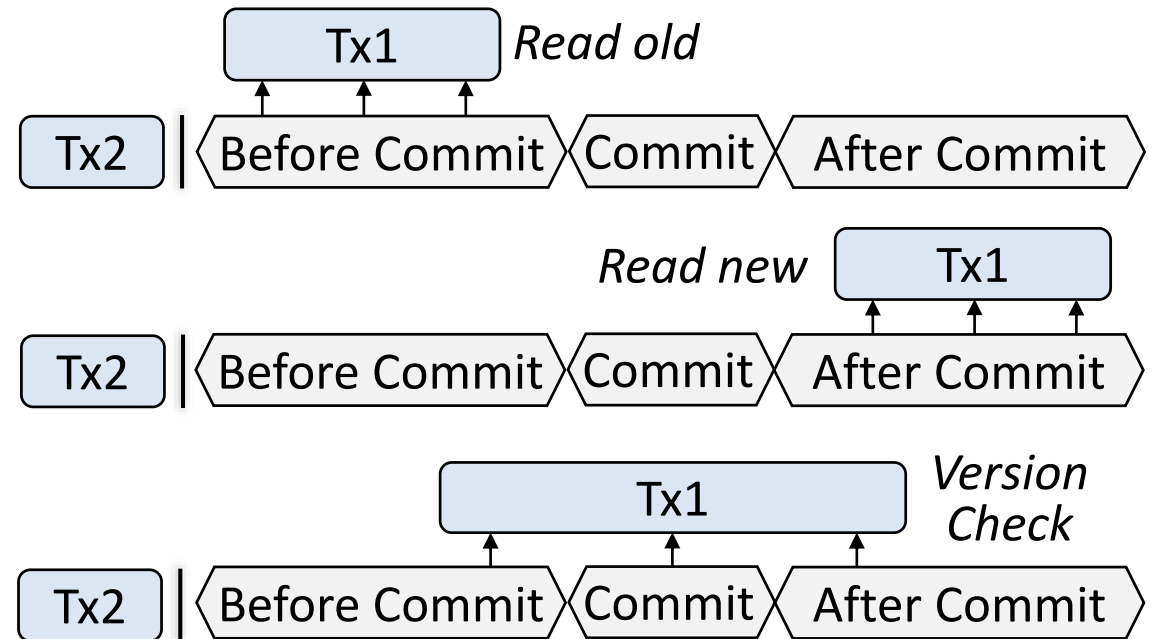
Backup-enabled Read

➤ Allows **backups** to serve **read-only (RO)** data

- In-place update → No address redirection
- Undo logging → No data migration → No CPU involvement in PM pool



Balance load → Improve throughput



Correctness guarantee

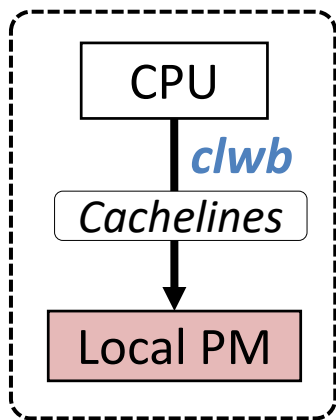
Selective Remote Flush

- Only issue remote flushes to backups after the final write
 - Guarantee remote persistency
 - Ensure recoverability by backups
 - Reduce flushing round trips
 - Compatible with different flush primitives ^[1]

¹One-sided RDMA **FLUSH** or **READ** (*READ-after-WRITE*)

Selective Remote Flush

- Only issue remote flushes to backups after the final write
 - Guarantee remote persistency
 - Ensure recoverability by backups
 - Reduce flushing round trips
 - Compatible with different flush primitives ^[1]

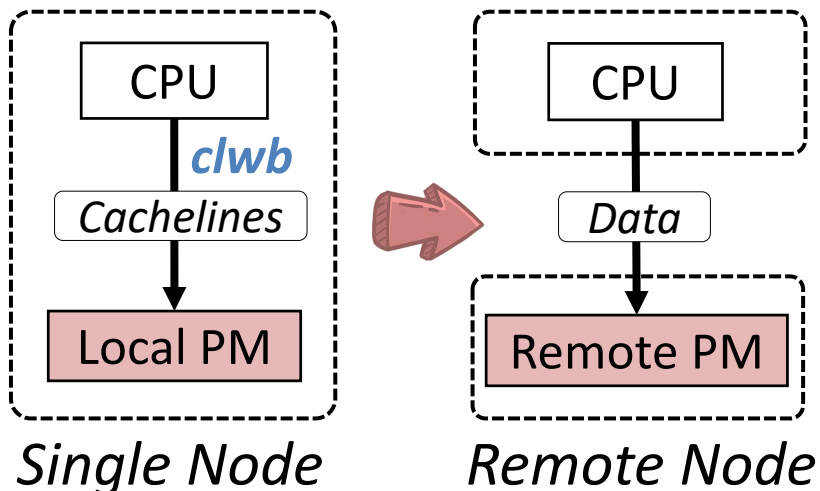


Single Node

¹ One-sided RDMA **FLUSH** or **READ** (*READ-after-WRITE*)

Selective Remote Flush

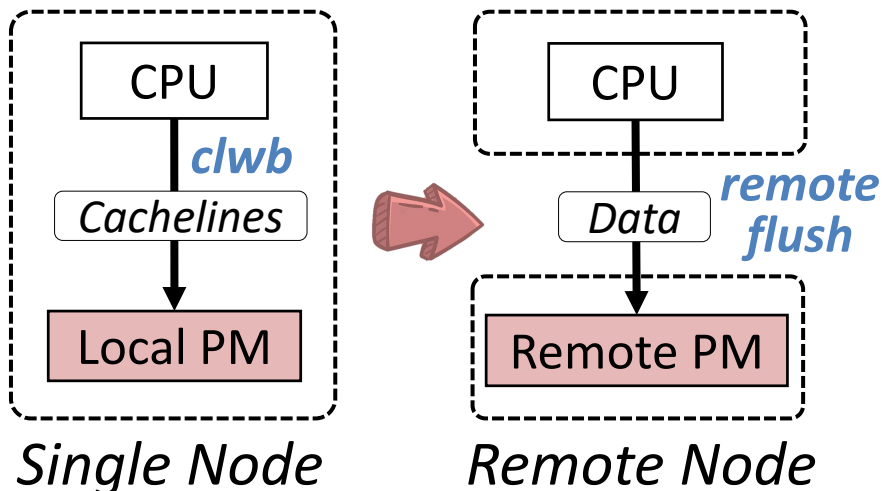
- Only issue remote flushes to backups after the final write
 - Guarantee remote persistency
 - Ensure recoverability by backups
 - Reduce flushing round trips
 - Compatible with different flush primitives ^[1]



¹One-sided RDMA **FLUSH** or **READ** (*READ-after-WRITE*)

Selective Remote Flush

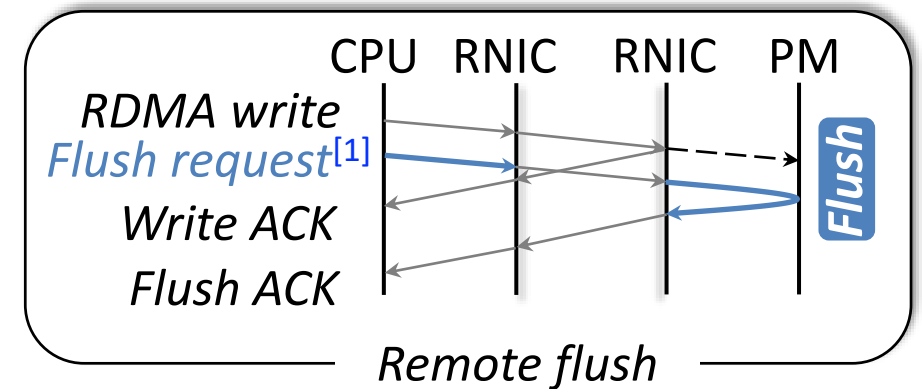
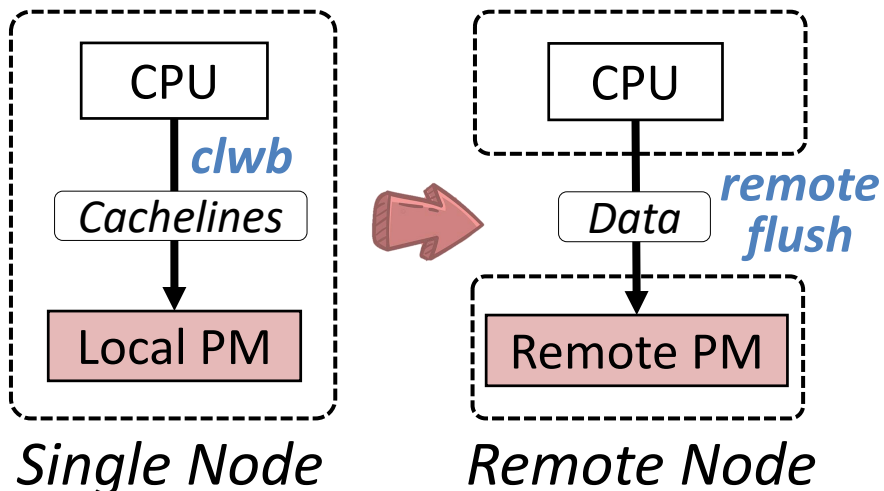
- Only issue remote flushes to backups after the final write
 - Guarantee remote persistency
 - Ensure recoverability by backups
 - Reduce flushing round trips
 - Compatible with different flush primitives ^[1]



¹ One-sided RDMA **FLUSH** or **READ** (*READ-after-WRITE*)

Selective Remote Flush

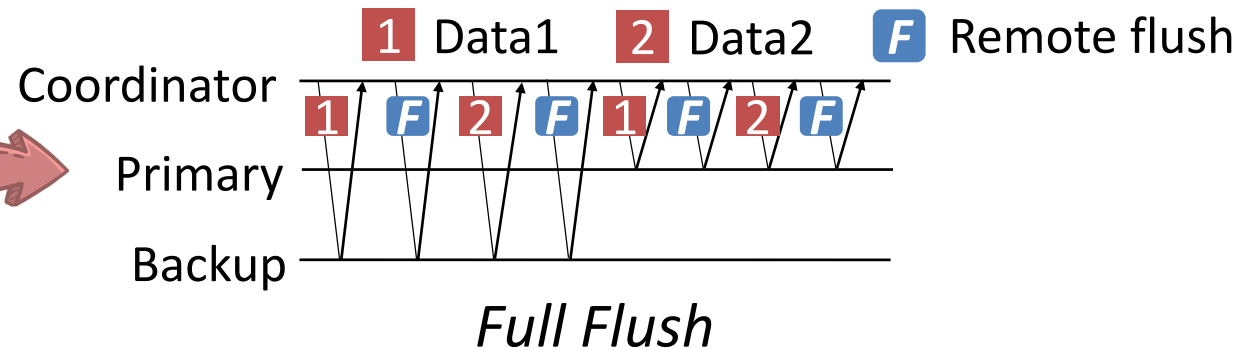
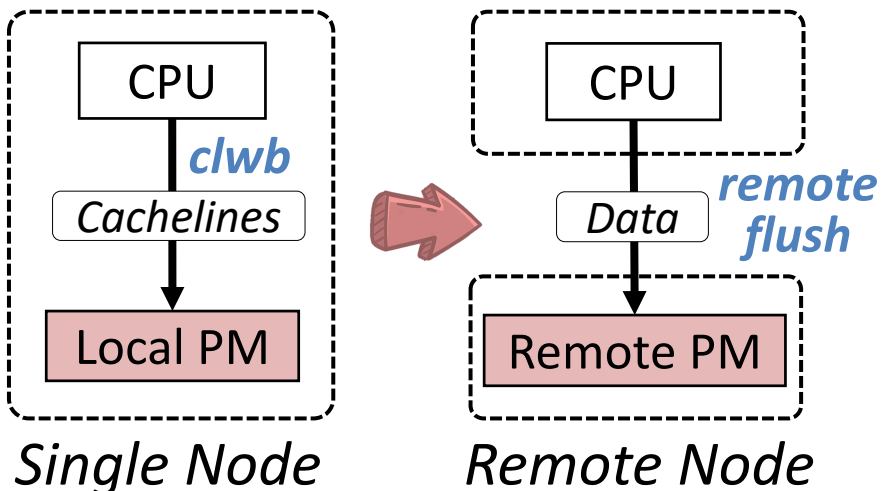
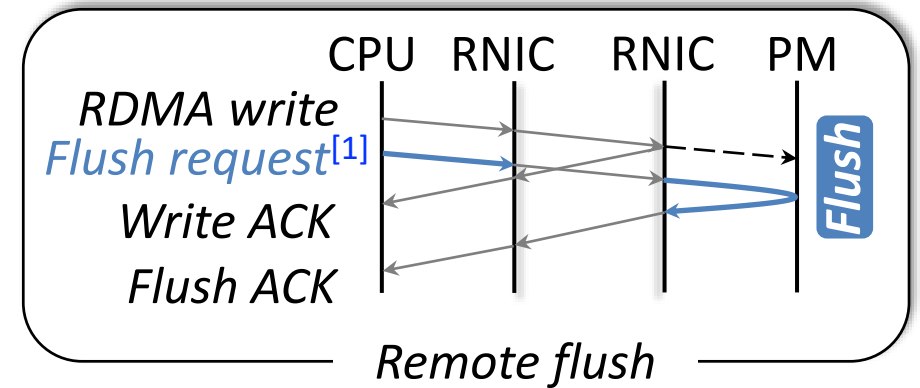
- Only issue remote flushes to backups after the final write
 - Guarantee remote persistency
 - Ensure recoverability by backups
 - Reduce flushing round trips
 - Compatible with different flush primitives ^[1]



¹One-sided RDMA **FLUSH** or **READ** (*READ-after-WRITE*)

Selective Remote Flush

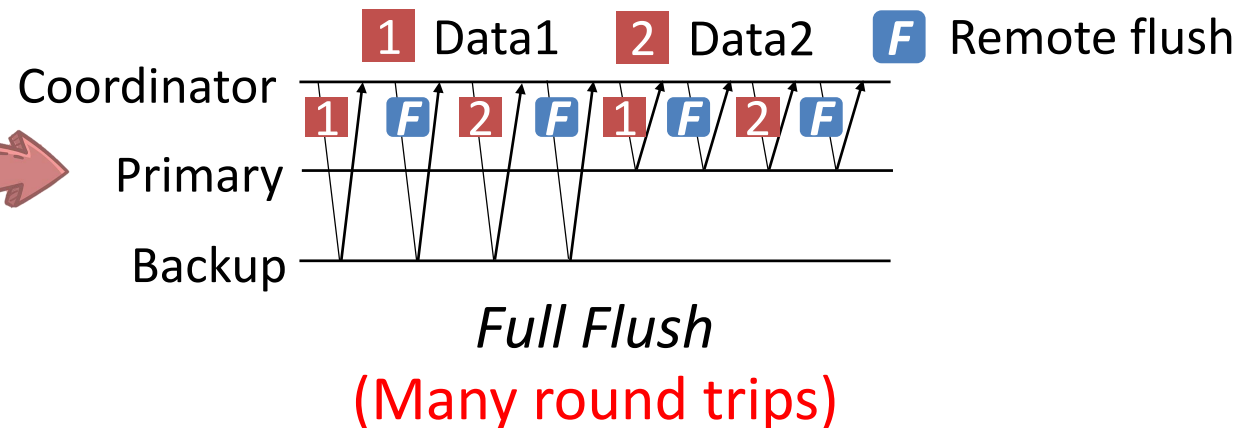
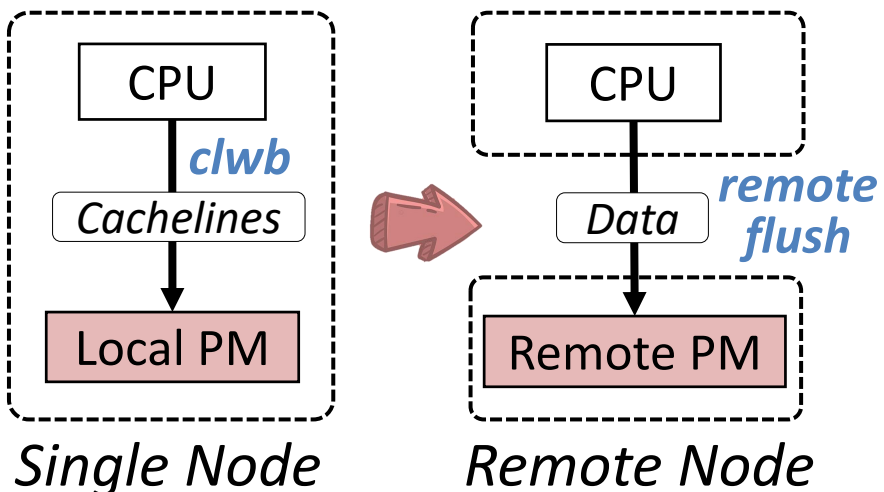
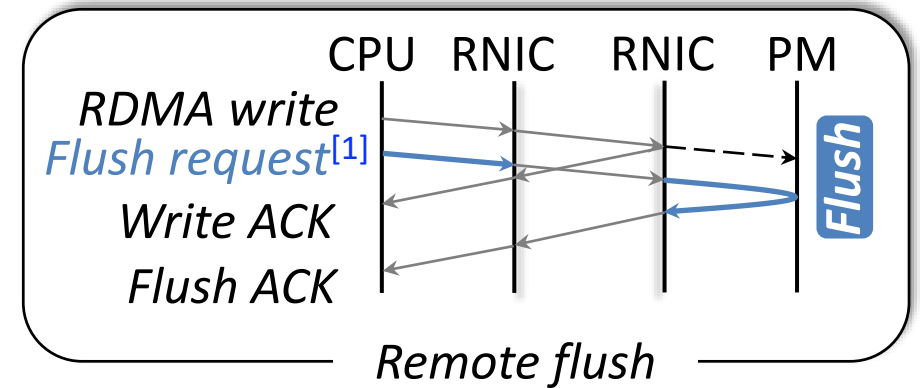
- Only issue remote flushes to backups after the final write
 - Guarantee remote persistency
 - Ensure recoverability by backups
 - Reduce flushing round trips
 - Compatible with different flush primitives ^[1]



¹ One-sided RDMA **FLUSH** or **READ** (*READ-after-WRITE*)

Selective Remote Flush

- Only issue remote flushes to backups after the final write
 - Guarantee remote persistency
 - Ensure recoverability by backups
 - Reduce flushing round trips
 - Compatible with different flush primitives ^[1]

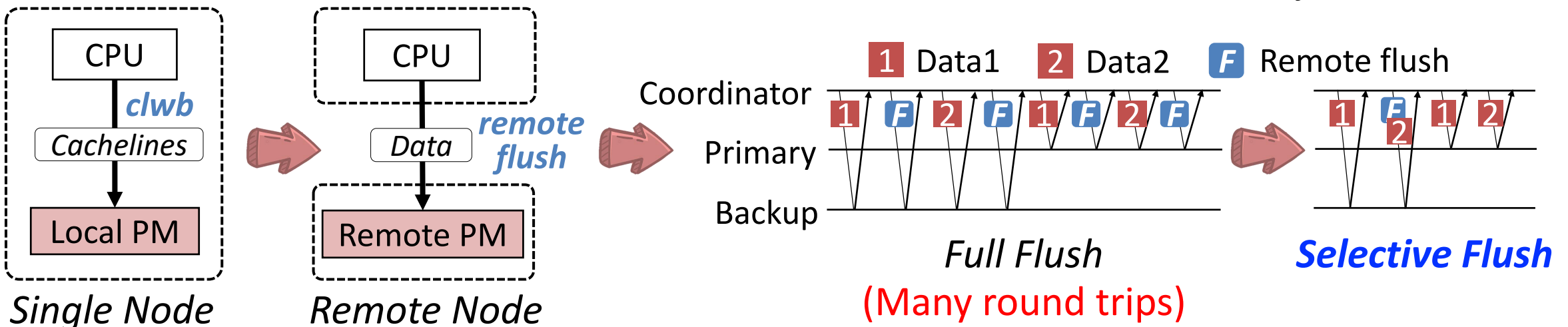
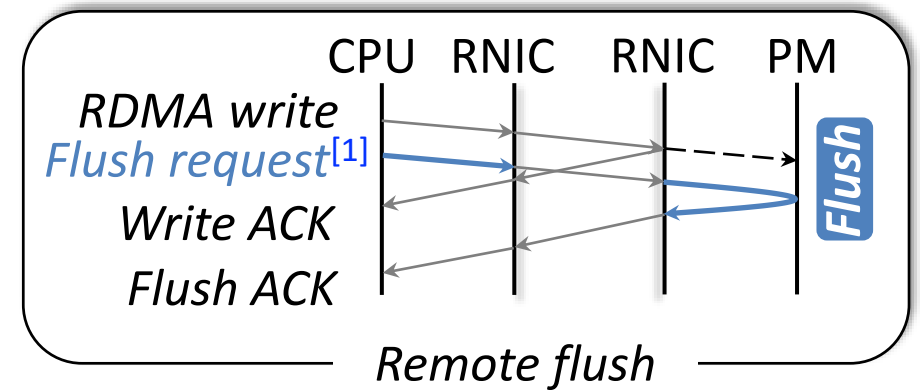


¹ One-sided RDMA **FLUSH** or **READ** (*READ-after-WRITE*)

Selective Remote Flush

➤ Only issue remote flushes to backups after the final write

- Guarantee remote persistency
- Ensure recoverability by backups
- Reduce flushing round trips
- Compatible with different flush primitives ^[1]

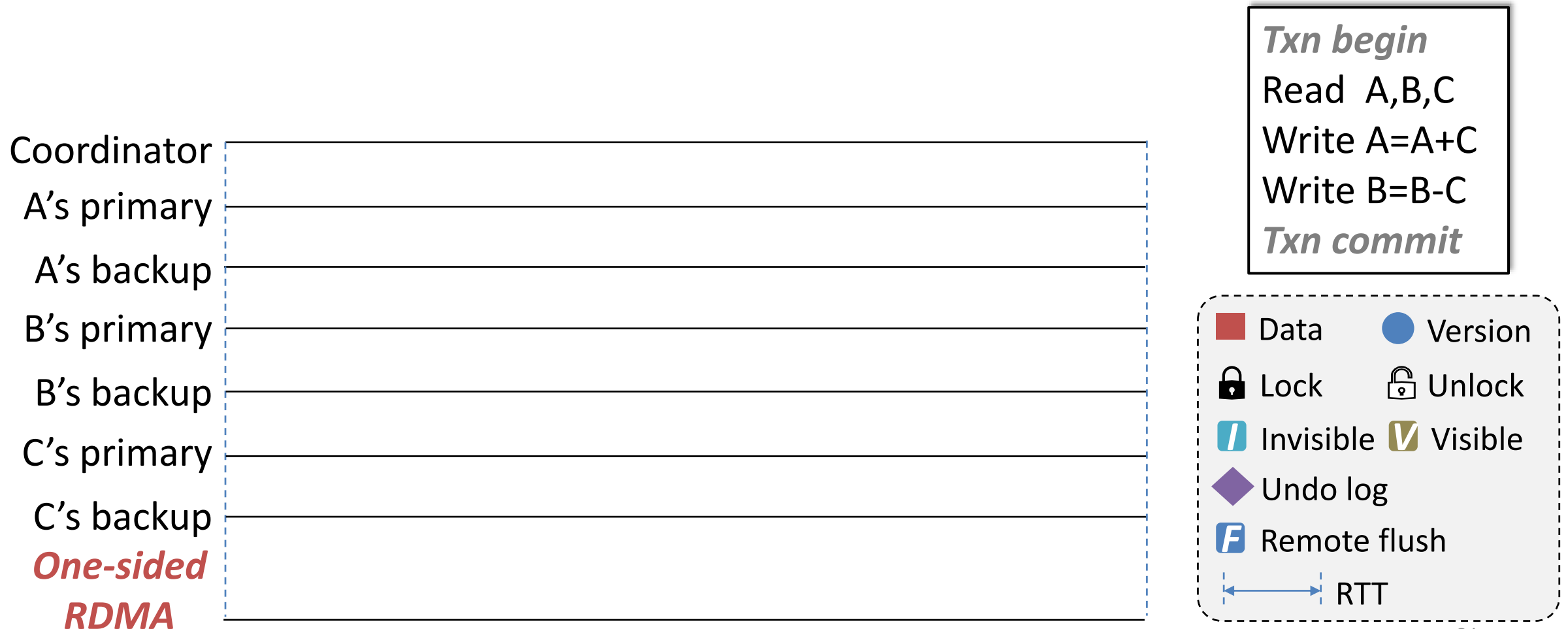


(Many round trips)

¹ One-sided RDMA **FLUSH** or **READ** (*READ-after-WRITE*)

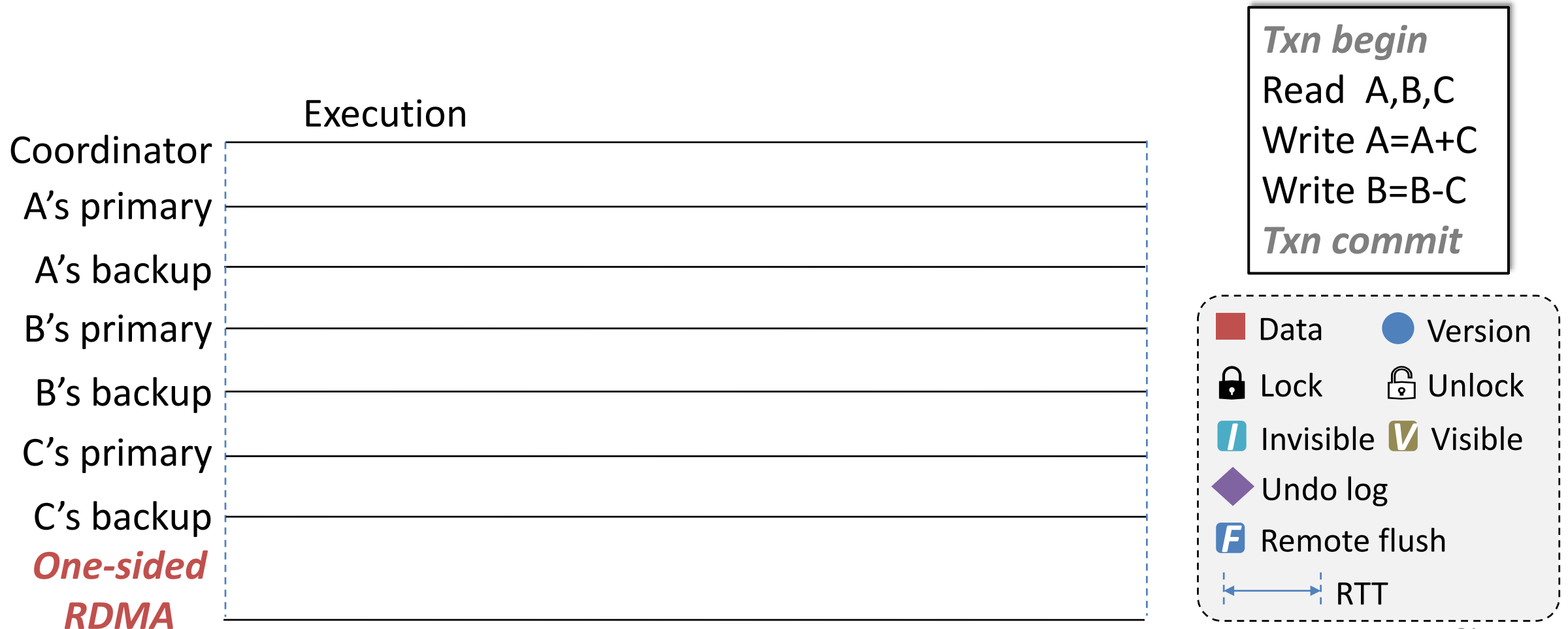
FORD Transaction

➤ **Put it all together:** One-sided RDMA and PM-conscious designs



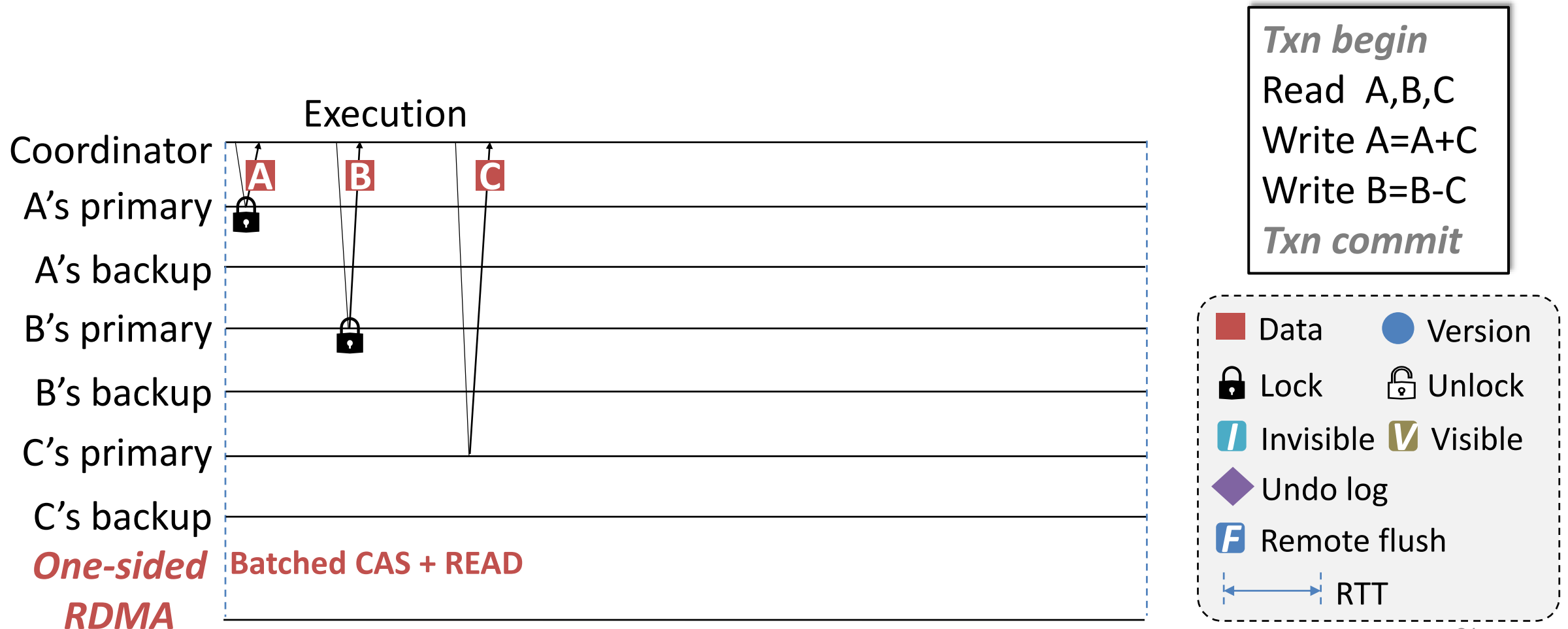
FORD Transaction

➤ **Put it all together:** One-sided RDMA and PM-conscious designs



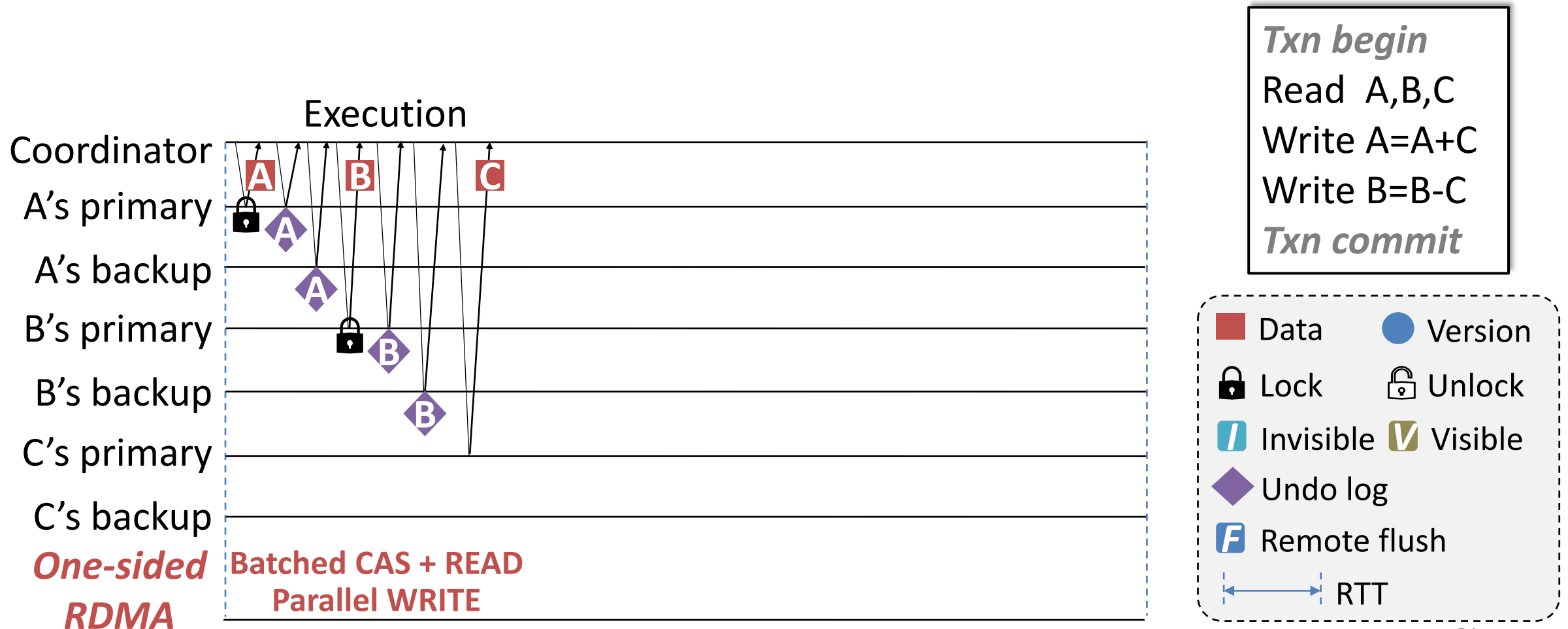
FORD Transaction

➤ **Put it all together:** One-sided RDMA and PM-conscious designs



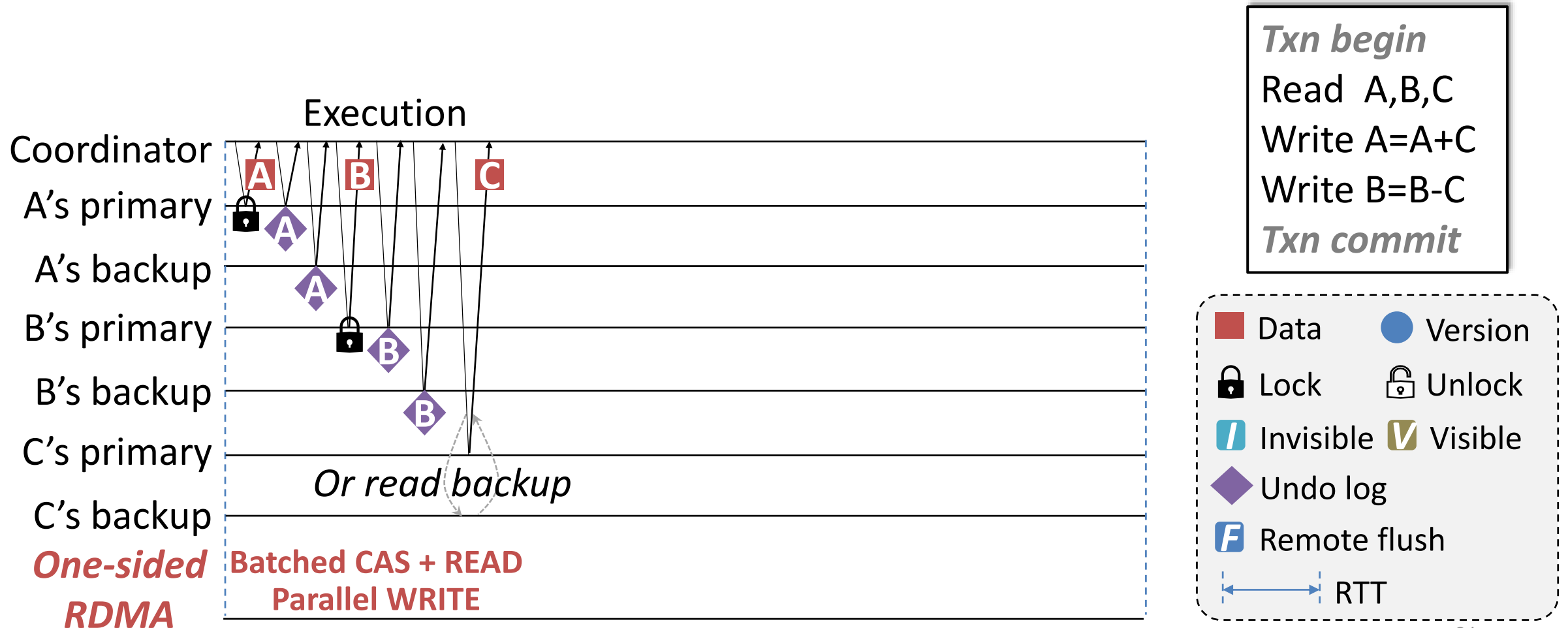
FORD Transaction

➤ **Put it all together:** One-sided RDMA and PM-conscious designs



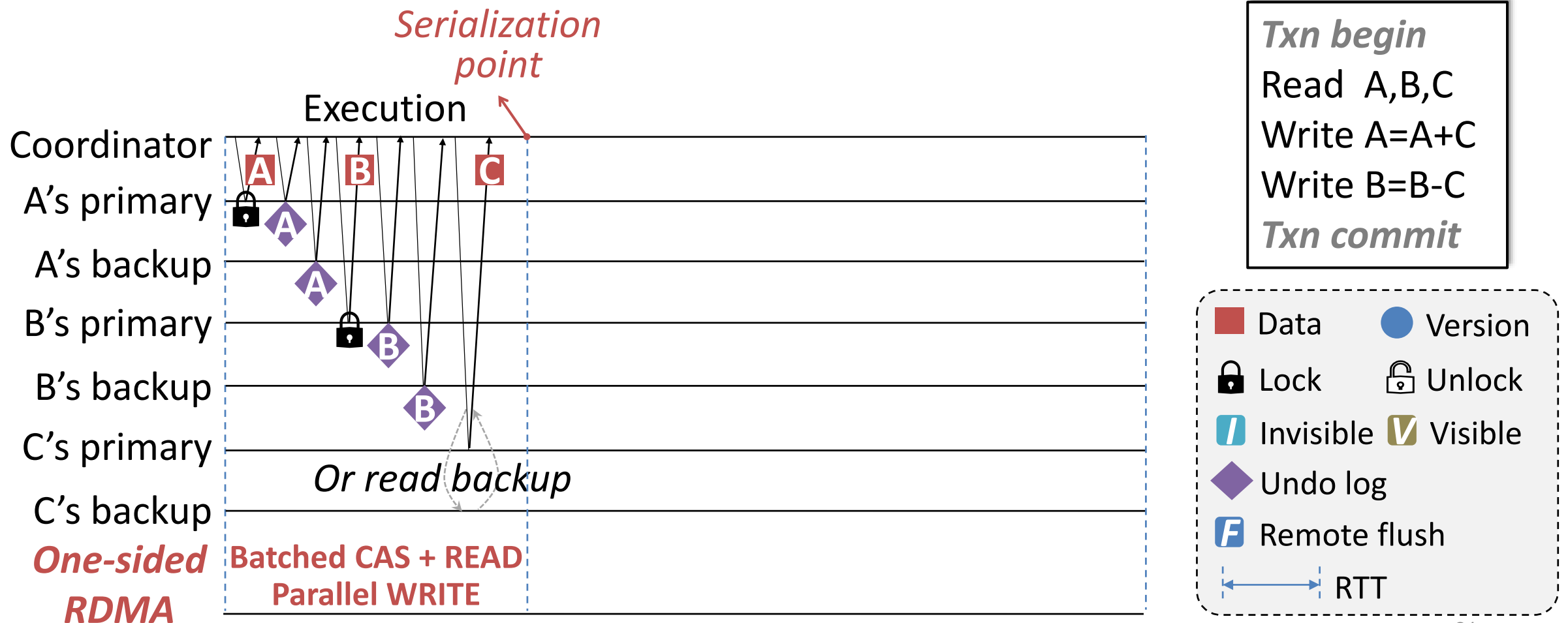
FORD Transaction

➤ **Put it all together:** One-sided RDMA and PM-conscious designs



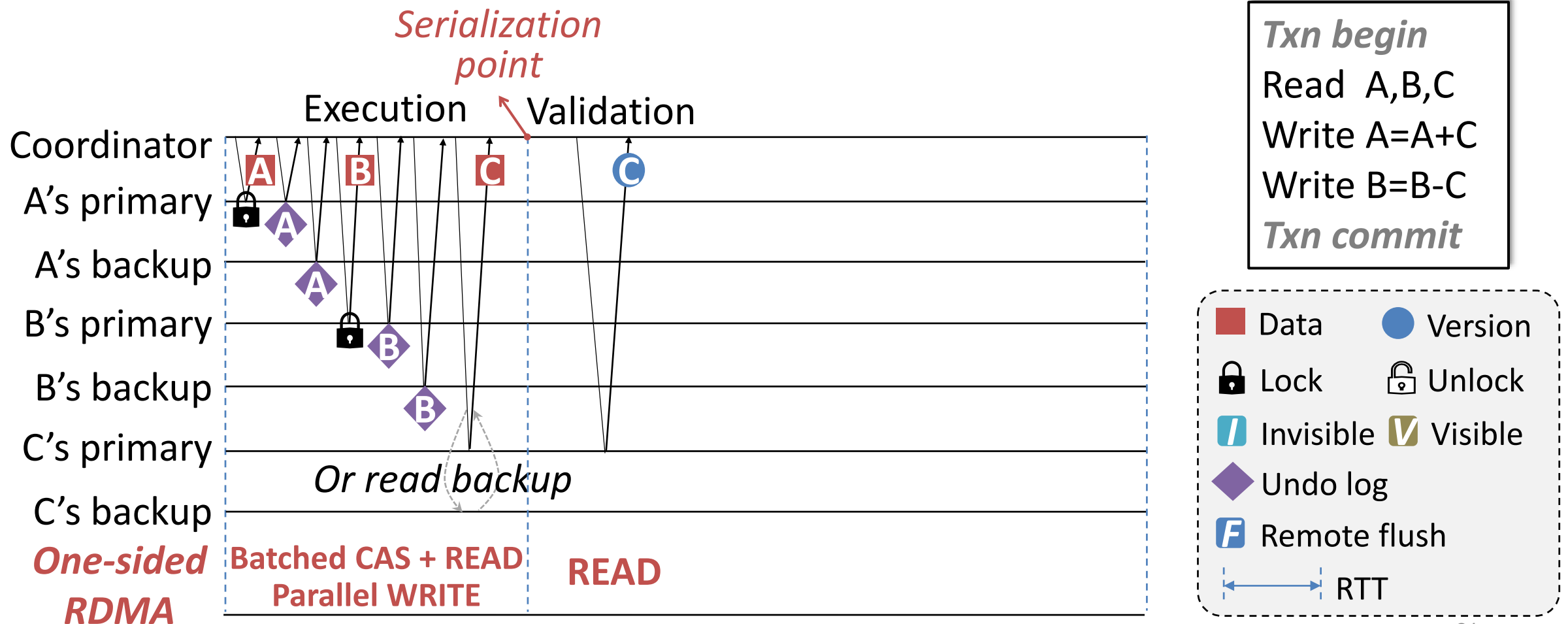
FORD Transaction

➤ **Put it all together:** One-sided RDMA and PM-conscious designs



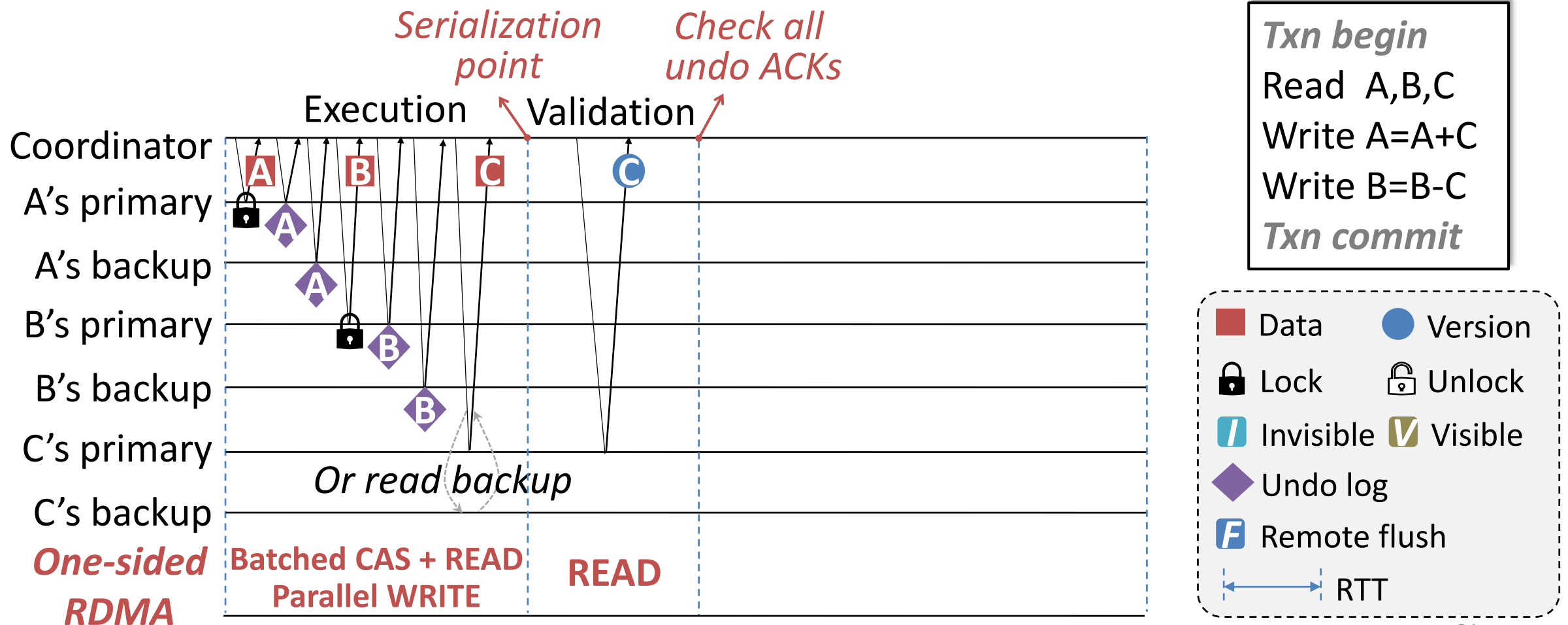
FORD Transaction

➤ **Put it all together:** One-sided RDMA and PM-conscious designs



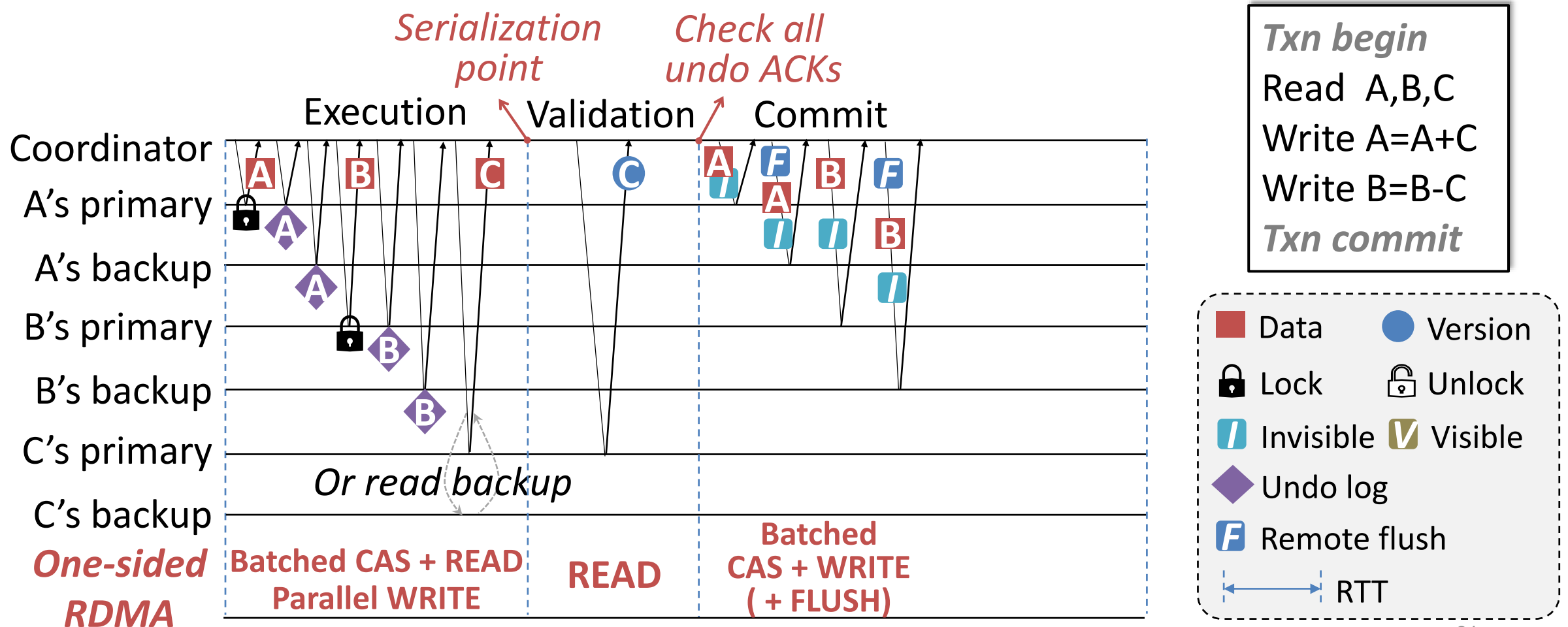
FORD Transaction

➤ **Put it all together:** One-sided RDMA and PM-conscious designs



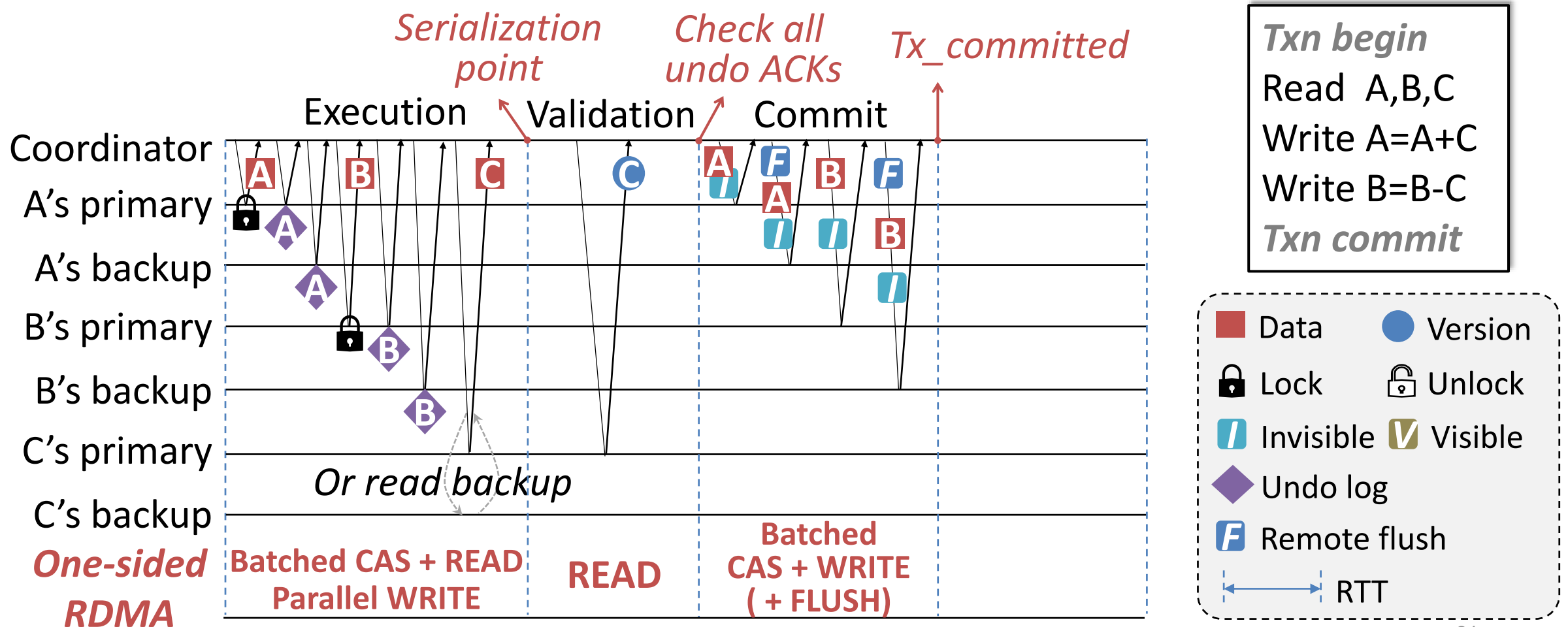
FORD Transaction

➤ **Put it all together:** One-sided RDMA and PM-conscious designs



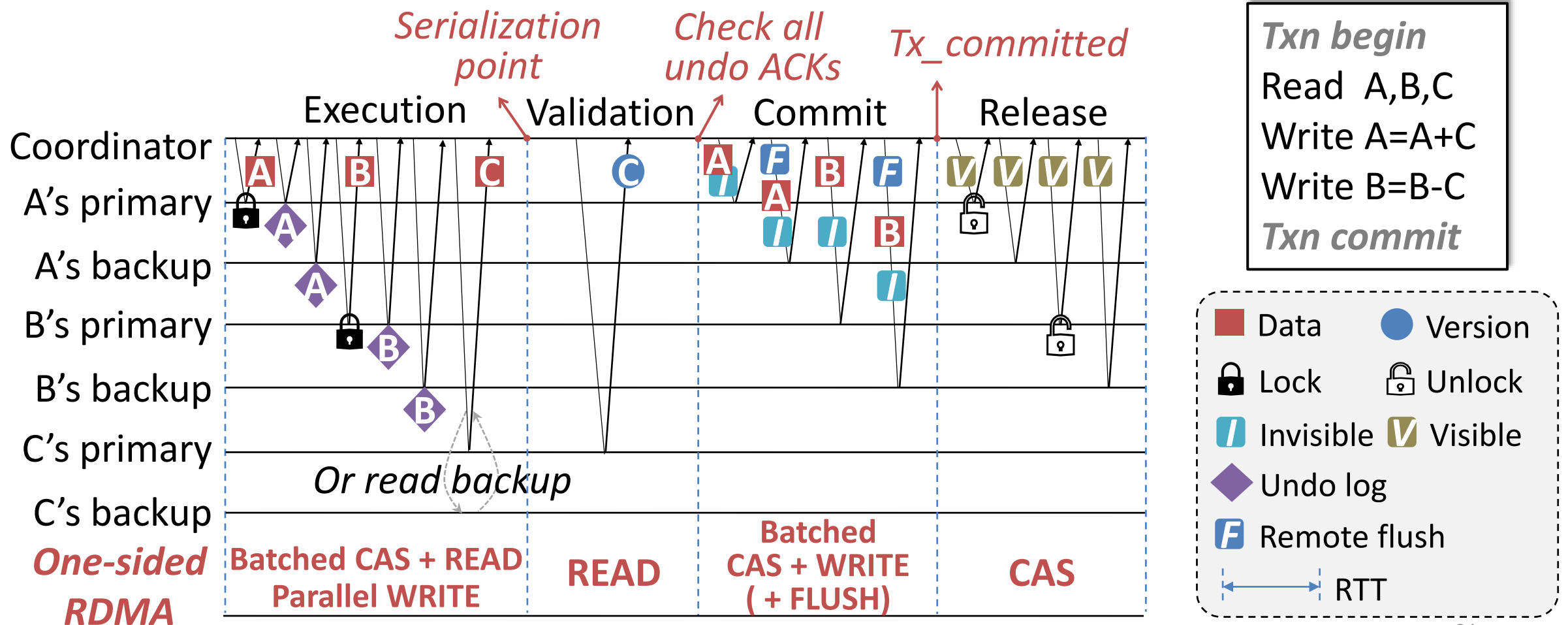
FORD Transaction

➤ **Put it all together:** One-sided RDMA and PM-conscious designs



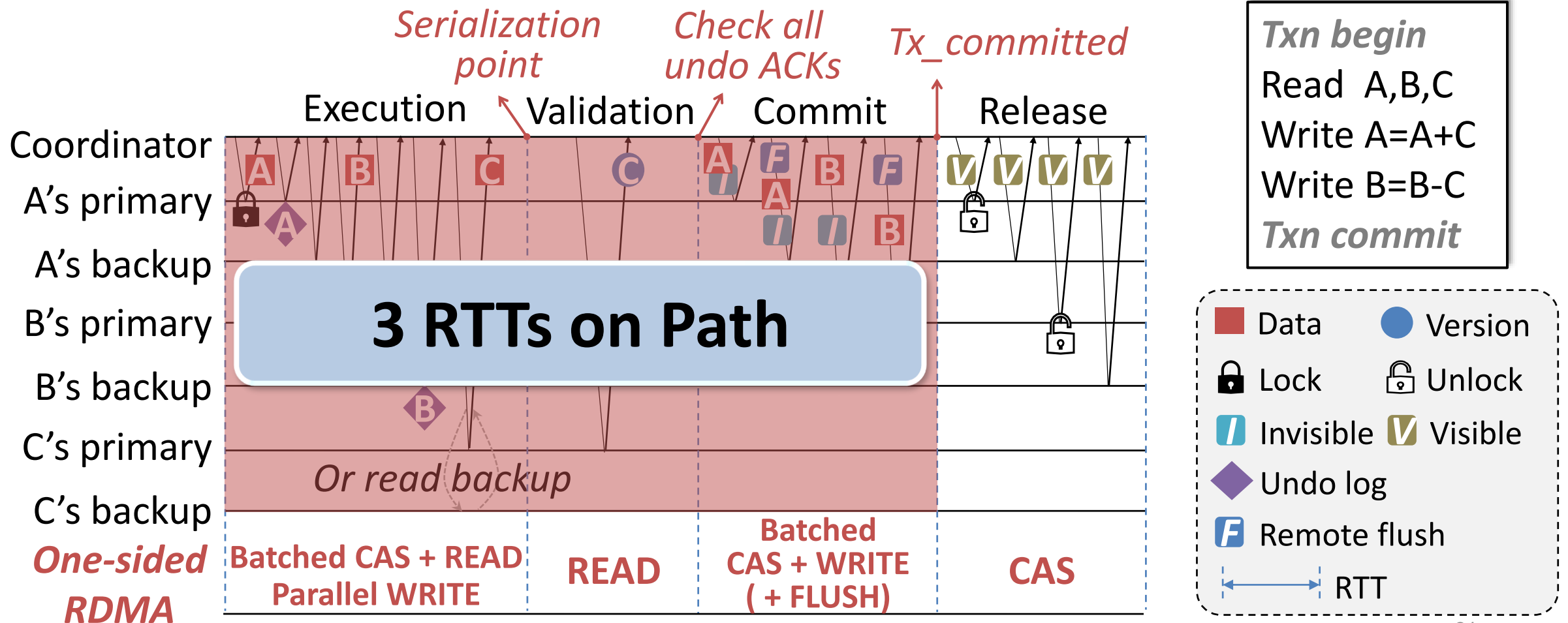
FORD Transaction

➤ **Put it all together:** One-sided RDMA and PM-conscious designs



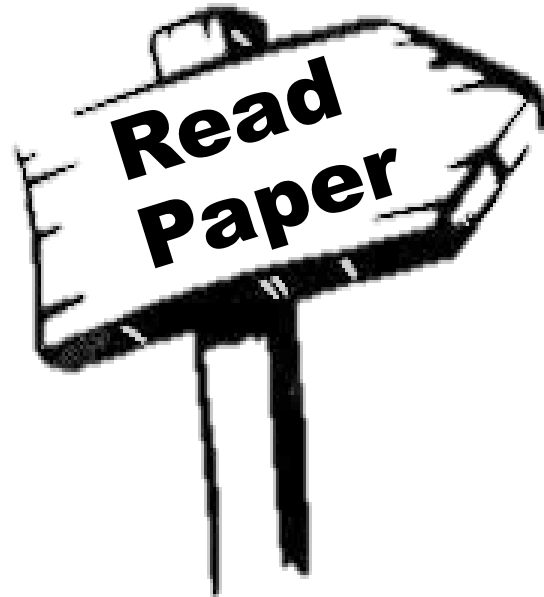
FORD Transaction

➤ **Put it all together:** One-sided RDMA and PM-conscious designs



More Details

- Programming Interface
- Indexes in PM Pool
- Analysis on ACID and Serializability
- ...

A collage of document pages. The top page shows a diagram of a hash table structure with buckets and slots. Below it is a list of steps (S1-S8) for calculating bucket addresses. The middle page contains C++ code for a WriteCheck transaction. The bottom page shows a comparison of sequential vs. interleaved transaction processing with a timeline diagram. A vertical text on the right edge of the collage reads: "phase, the actor clears not be re- 1)-way", "ing undo", "ures by", "ORD", "e.g.,", "stoc-", "ness", "all".

RDMA-registered memory region

Bucket 1 Bucket 2 ... Bucket N Bucket 1

Slot FObj1 FObj2 ... FObjN

Figure 8: The hash table structure in FStore.

To enable coordinators to calculate remote addresses for one-sided RDMA in dtxn processing, the connection manager in PM pool sends the metadata (as listed below) of each hash table to the compute pool during network interconnections.

- TableID (8B): Global unique database table id.
- Addr (8B): Virtual start address of this hash table.
- Off (8B): Offset between Addr and MR's start address.
- BucketNum (8B): Bucket number of the hash table.
- BucketSize (4B): Size of a bucket (in bytes).
- SlotNum (4B): Number of slots per bucket.

Given the key (e.g., K0) of a record, if its remote address is buffered in the local cache, the coordinator directly reads the record using an RDMA READ. Otherwise, the coordinator reads a remote record as the following Steps:

S1: Calculate the bucket id:
 $bucket_id = Hash(K0) \bmod BucketNum$

S2: Calculate the bucket offset in the remote MR:
 $bucket_off = bucket_id \times BucketSize + Off$

S3: Read the remote bucket (*bkr*) using *bucket_off*.

S4: Compare *K0* with the *SlotNum* keys in *bkr*. If a key = *K0*, the record is obtained. Then go to S7. Or else go to S5.

S5: If the next field of *bkr* is NULL, there is no such remote record. Then go to S8. Or else go to S6.

S6: Calculate the next bucket offset as below and go to S3.
 $bucket_off = bkr_next - Addr + Off$

S7: Exit if the record is visible. Or else re-read it until visible.

S8: Exit with a *KEY_NOT_EXIST* hint.

Since the metadata size of a hash table is only 40B and each remote address is 8B, the local cache in compute pool can buffer all these metadata and addresses, as shown in Fig. 15b. Caching metadata is scalable, because the compute blades do not need to synchronize their metadata with each other: 1) The metadata of index does not change. 2) If the cached addresses are stale, FORD enables the coordinator to detect this and update its own cached addresses, as discussed in § 3.2.

4.2 Transaction Interfaces

FORD provides a runtime library, called *FLib*, for applications to process dtxns. Flib exposes the following interfaces:

- TxBegin: Start to execute a dtxn and record its id.
- AddrRO: Add an initialized FObj to the read-only set.
- AddrRW: Add an initialized FObj to the read-write set.
- TxExecute: The coordinator reads the remote data specified in read-only and read-write sets, and then executes the dtxn logic. Our hitchhiked locking and backup-enabled read schemes are leveraged.
- TxCommit: After execution, the coordinator commits the updated data to remote primaries and backups using our coalescent commit and selective remote flush schemes.

```
1 bool WriteCheck(uint64_t dtxn_id, DTXID* dtxn) {
2 // The dtxn involves FLib interfaces
3 dtxn->TxBegin(dtxn_id);
4 // Use a random account as the key
5 uint64_t acct_id = RandomAccount();
6 FObj* sav_obj = new FObj(SavingsTableID, acct_id);
7 FObj* chk_obj = new FObj(CheckingTableID, acct_id);
8 dtxn->Addr(sav_obj);
9 dtxn->Addr(chk_obj);
10 if (!dtxn->TxExecute()) return false;
11 // Get record values and run transaction logic
12 sav_val_t* sav = (sav_val_t*) sav_obj->value;
13 chk_val_t* chk = (chk_val_t*) chk_obj->value;
14 if (sav->balance + chk->balance < PredefineAmount)
15   chk->balance += PredefineAmount + 1;
16 else chk->balance -= PredefineAmount;
17 bool status = dtxn->Commit();
18 delete sav_obj; delete chk_obj;
19 // Report commit (true) or abort (false) to client
20 return status;
21 }
```

Figure 9: The example of C++ code using FLib interfaces.

The CPU runtime in the Execution, Validation, and Commit phases.

Wait ACK / Execution flow

Time

Tx1 Tx2 Tx3

(a) Sequential processing

(b) Interleaved processing

Figure 10: The comparisons between (a) sequential processing, and (b) interleaved processing, in one thread.

Our transaction interfaces support general transaction processing. Specifically, the developers are not required to know all the read/write sets at the beginning of each transaction. Instead, developers call *AddrRO*, *AddrRW*, and *TxExecute* multiple times when reading/writing data occurs during a transaction. Fig. 9 illustrates an example of using our interfaces in the Write Check transaction of the SmallBank benchmark [50]. This transaction reads the balances from the Savings and Checking tables, and updates the balance in the Checking table. It shows that our interfaces are easy-to-use.

4.3 Interleaved Transaction Processing

As shown in Fig. 10a, sequentially processing dtxns in a thread wastes the CPU cycles due to waiting for RDMA ACKs, which significantly decreases the throughput. To avoid CPU idling in the compute pool, FORD leverages an interleaved processing scheme that enables multiple coordinators in one thread to process different dtxns in pipeline, as presented in Fig. 10b. In this way, the network RTTs are efficiently hidden and the CPU cores in the compute pool are fully utilized to improve the throughput.

We use coroutines [29, 60] to implement the interleaved processing. Each CPU thread generates several coroutines and each coroutine acts as a coordinator to execute dtxns. After issuing the RDMA requests, a coroutine yields its CPU core to another coroutine to process the next dtxn. A dedicated coroutine in each thread polls RDMA ACKs. If all ACKs of a coroutine arrived, FORD schedules this coroutine to occupy the CPU core to resume execution. The results in Fig. 16 show that using a proper number of coroutines improves the throughput without heavily increasing the latency.

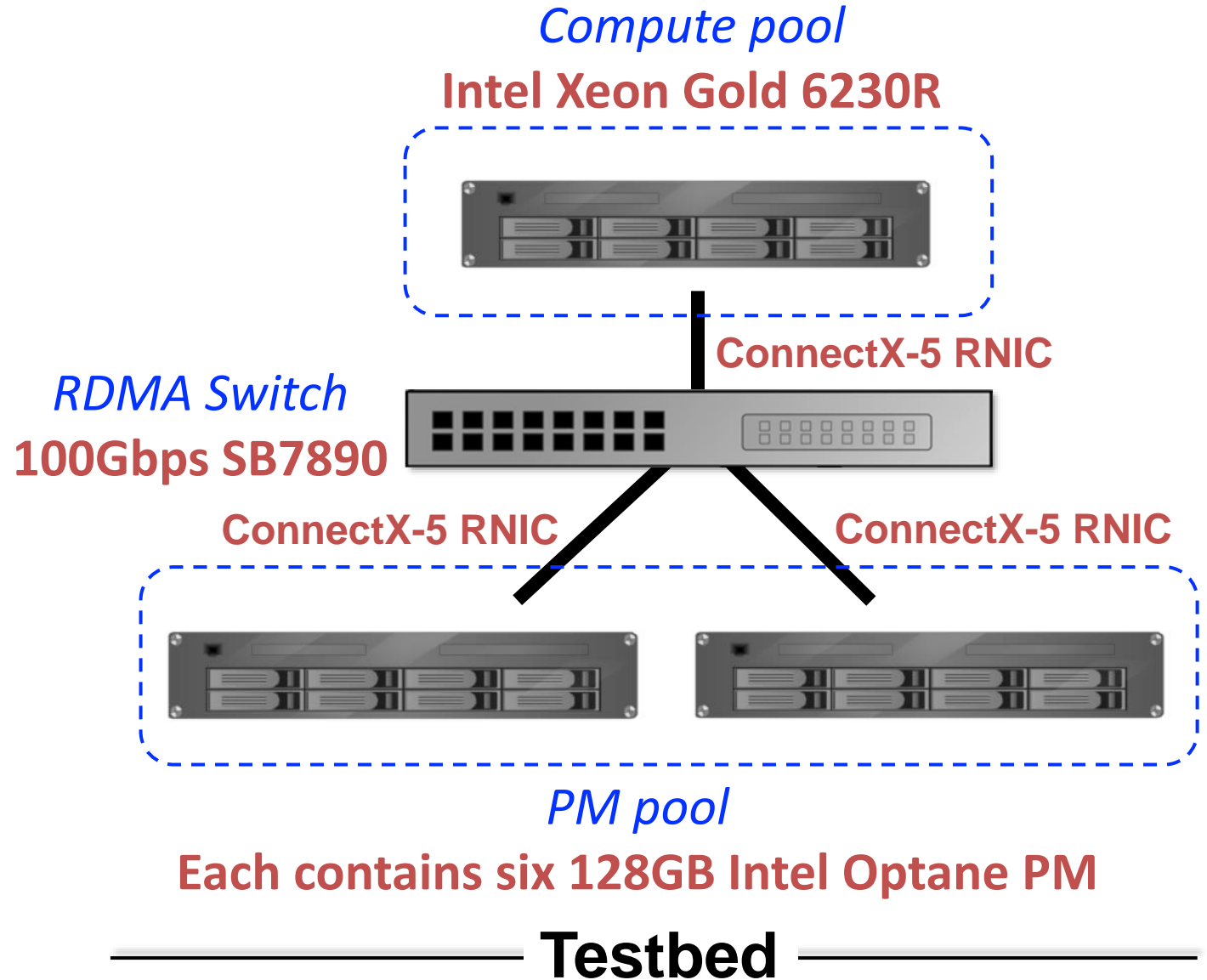
Evaluation

➤ Benchmarks

- KV store
 - 8B key + 40B value
 - Skewed ($\theta = 0.99$) + Uniform
- TATP
 - RO/RW: 80%/20%, 48B
- SmallBank
 - RO/RW: 15%/85%, 16B
- TPCC
 - RO/RW: 8%/92%, 672B

➤ Comparisons^[1]

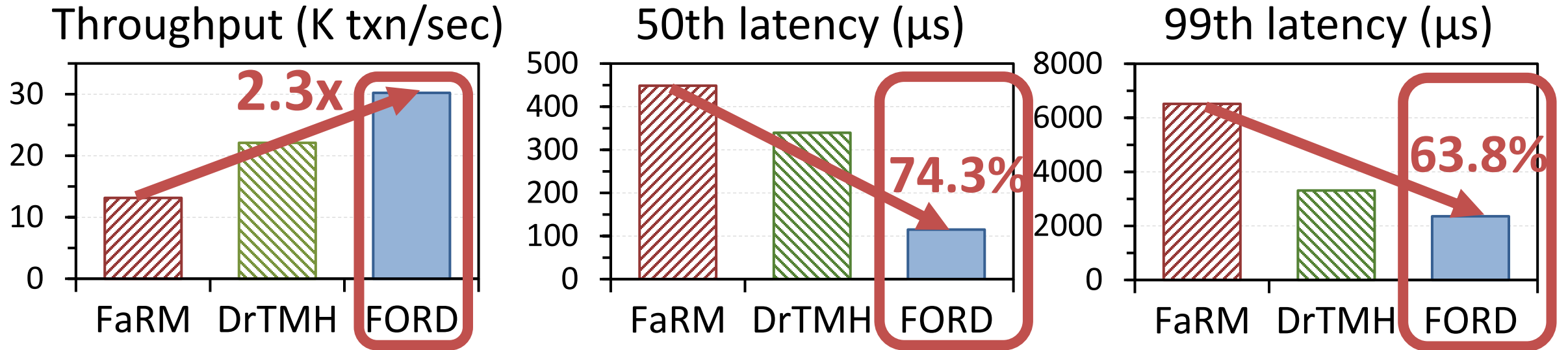
- FaRM@SOSP'15
- DrTM+H@OSDI'18



¹ Protocols are re-implemented using one-sided RDMA

End-to-End Performance

- 112 coordinators
- **Efficient round trip reduction and backup utilization**

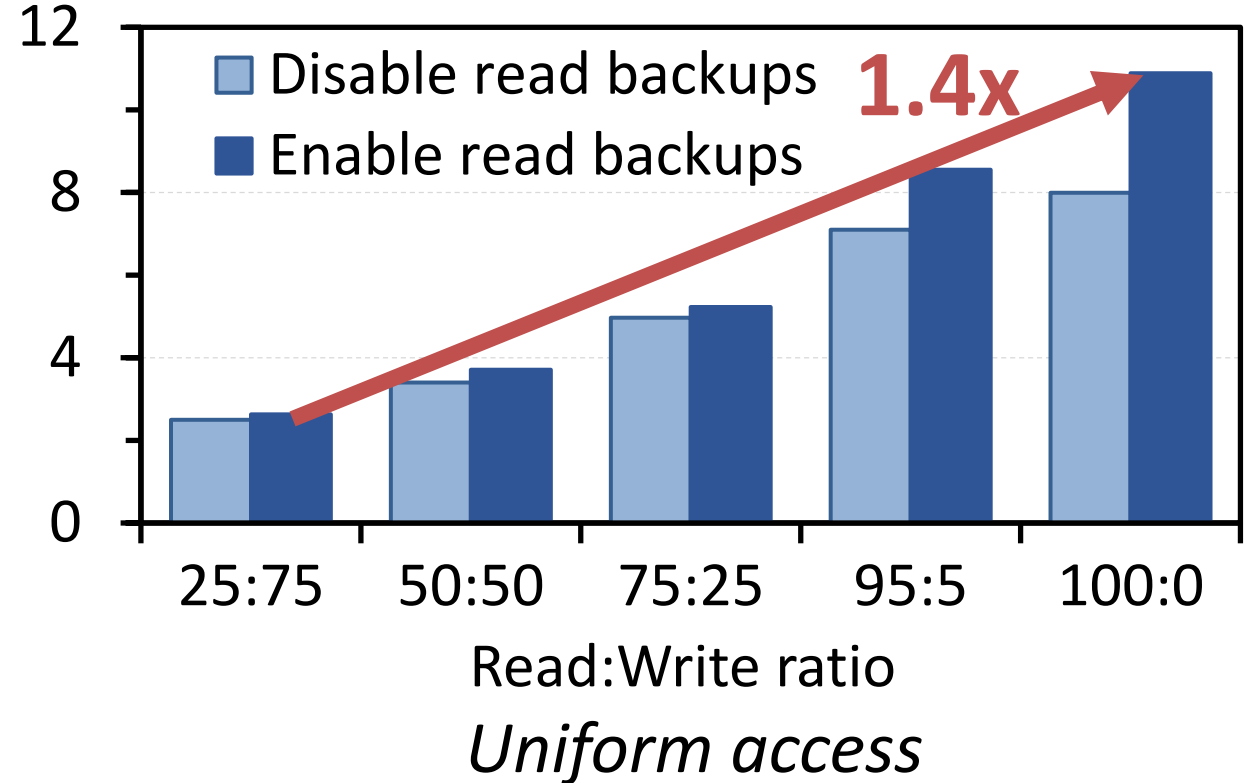
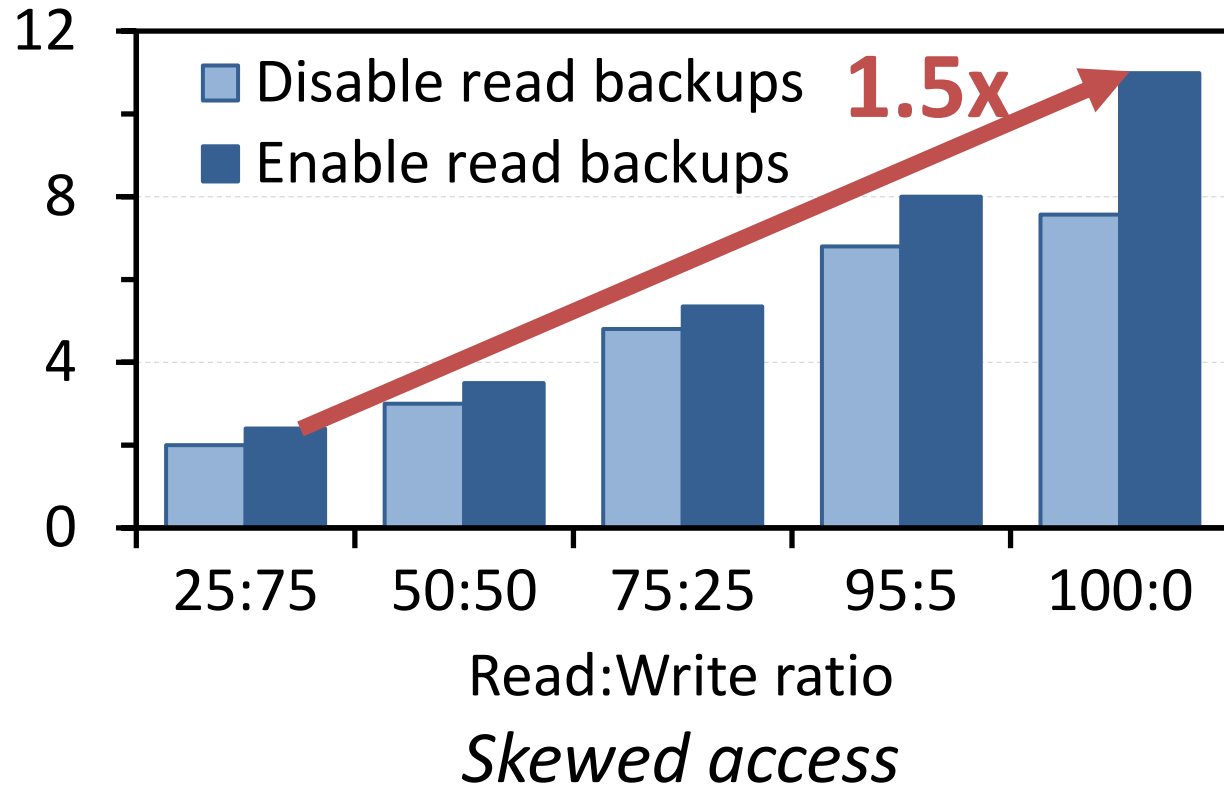


TPCC results

Read backup

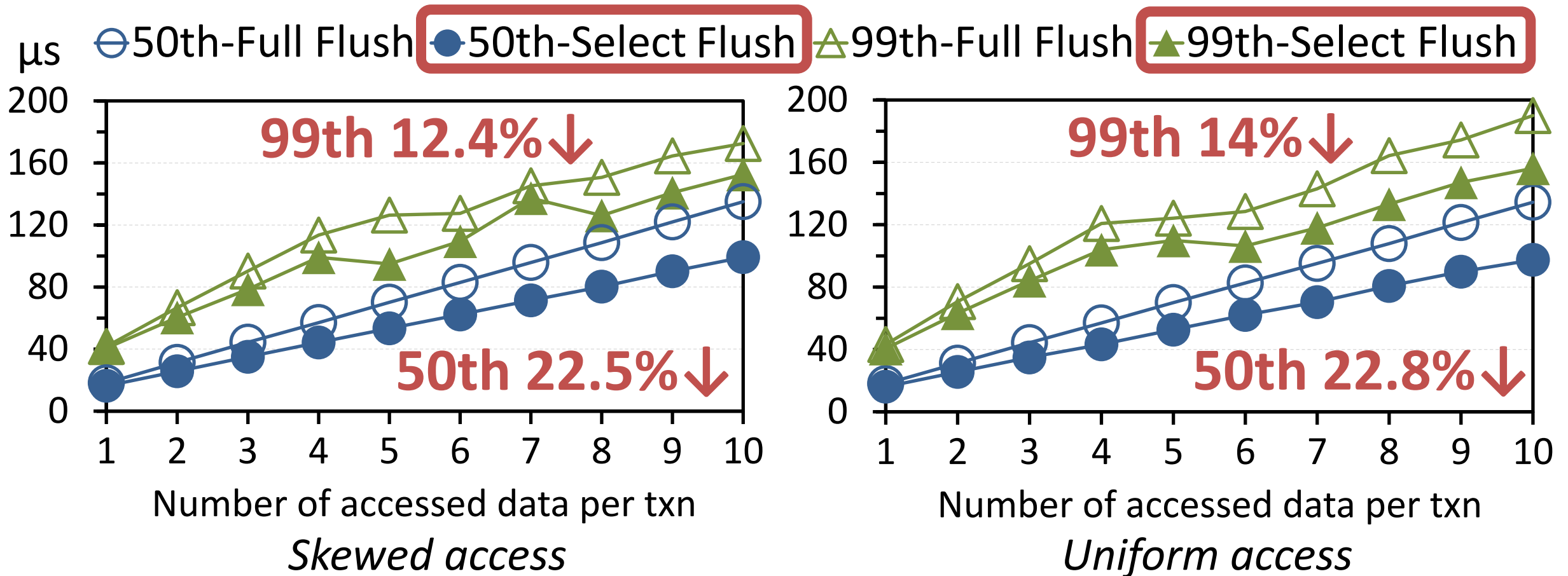
➤ KV store, 1 backup

M txn/sec

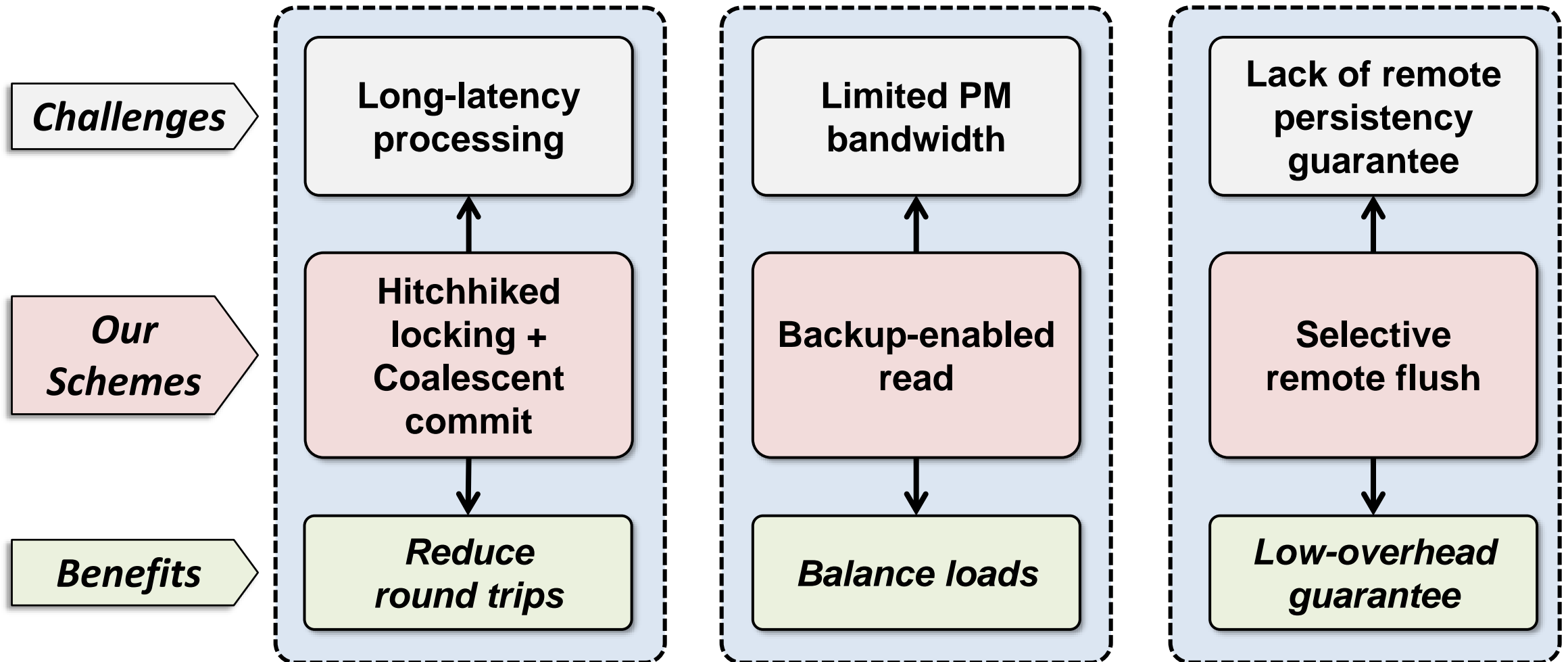


Remote flush

- KV store, 1 coordinator



Conclusion



Thank you! Q&A